1
2
3
# Verification of TIGGE Multi-model and ECMWF Reforecast-Calibrated Probabilistic Precipitation Forecasts over the Contiguous US

6
7

Thomas M. Hamill

*NOAA Earth System Research Laboratory, Physical Sciences Division*

*Boulder, Colorado USA*

13
14

21
22

24
25
26
27
28
29
30
31
32
33
34

Corresponding author address:

Dr. Thomas M. Hamill
NOAA ESRL, Physical Sciences Division
R/PSD 1
325 Broadway
Boulder, CO 80305-3328
tom.hamill@noaa.gov
phone: (303) 497-3060  fax: (303) 497-6449

44
45                                               ABSTRACT
46
47          Probabilistic quantitative precipitation forecasts (PQPFs) were generated

48    during July to October 2010 using European Centre (ECMWF), United Kingdom

49    (UKMO), US (NCEP) and Canadian (CMC) forecast data.  24-hour accumulated

50    precipitation forecasts were evaluated within the contiguous US against

51    precipitation analyses for +1 to +5 days lead at 1-degree grid spacing.

52          PQPFs from ECMWF's ensembles generally had the highest skill of the raw

53    ensemble forecasts, followed by CMC.  PQPFs from CMC forecasts were the most

54    reliable but the least sharp. PQPFs from NCEP and UKMO ensembles were the least

55    reliable but sharper.

56          Multi-model PQPFs were more reliable and skillful than individual ensemble

57    prediction system forecasts.  The improvement was larger for heavier precipitation

58    events than light events.

59          ECMWF ensembles were statistically post-processed using extended logistic

60    regression and the five-member weekly reforecasts for the June - November period

61    of 2002-2009.  Multi-model ensembles were also post-processed using logistic

62    regression and the last 30 days of prior forecasts and analyses. The reforecast-

63    calibrated ECMWF PQPFs were more skillful and reliable for the heavier

64    precipitation events than ECMWF raw forecasts but less sharp.  Raw multi-model

65    PQPFs were generally more skillful than reforecast-calibrated ECMWF PQPFs for the

66    light precipitation events but about the same skill for the heavier events; also, they

67    were sharper but somewhat less reliable than ECMWF reforecast-based PQPFs.

68    Post-processed multi-model PQPFs did not provide as much improvement to the

69    raw multi-model PQPF as the reforecast-based processing did to the ECMWF

70    forecast.

71         The evidence presented here suggests that all operational centers, even

72    ECMWF, would benefit from generating reforecasts and sharing data in real time.

73

74

1. **Introduction**

76
77      An ongoing challenge with short- and medium range ensemble prediction

systems (EPSs) is how to generate probabilistic forecasts that account for the

system errors in the ensemble.   System errors include sampling error due to the

finite ensemble size, the error introduced by model imperfections such as the grid

truncation, the use of deterministic parameterizations (Houtekamer and Mitchell

2005), and assimilation system and observation imperfections.  There are many

methods for treating system error, from introducing stochastic aspects into the

ensemble prediction system (Buizza et al. 1999, Shutts 2005, Berner et al. 2009,

Palmer et al. 2009, Charron et al. 2010), using multiple parameterizations  (Charron

et al. 2010, Berner et al. 2011), using multiple models (Bougeault et al. 2010), and

statistical post-processing.

      Two methods that will be explored and contrasted here are the multi-model

methods and statistical post-processing.   The underlying hypothesis of multi-model

ensembles (Krishnamurti et al. 2000, Wandishin et al. 2001, Mylne et al. 2002,

Doblas-Reyes et al. 2005, Hagedorn et al. 2005, Weigel et al. 2008, Candille 2009,

Johnson and Swinbank 2009, Bougeault et al. 2010, Iversen et al. 2011) is that the

many differences between constituent EPSs will result in them generating ensemble

forecasts with quasi-independent systematic errors, so the combination may result

in a more accurate estimate of the uncertainty.   Practically, also, these are

ensembles of opportunity.  If all centers are willing to share rather than sell their

forecast data, the additional members can be used for only the cost of data

98    transmittal and storage, so they may provide an inexpensive way to improve

99    forecast skill.  However, there are some potential disadvantages of multi-model

100   ensembles.   Developing an accurate, stable weather prediction system is costly, so

101   multi-model ensembles are likely to be less useful when formed from immature

102   systems.  System outages may prevent routine access to other centers' ensembles.

103   One or other of the models is likely to have been changed recently, rendering it

104   difficult to understand the multi-model system error characteristics.  Also, the

105   hypothesis of quasi-independent errors may not always hold.   Practically, each

106   operational center is interested in providing a high-quality product without

107   depending on another center's data.  When another center develops a method that

108   improves the forecast significantly, it may be adopted at other operational centers.

109   The similarity could result in some co-linearity of errors and decreased collective

110   usefulness (Lorenz et al. 2011).

111        Another method for addressing system error is through statistical post-

112   processing.  Discrepancies between time series of past forecasts from a fixed model

113   and the verifying observations/analyses can be used to modify the real-time

114   forecasts.   For some variables such as short-range forecasts of surface temperature,

115   a short time series may be sufficient (Stensrud and Yussouf 2003, Yussouf and

116   Stensrud 2007, Hagedorn et al. 2008).  For others such as heavy precipitation and

117   longer-lead forecasts, using a long time series of reforecasts has been shown to

118   dramatically improve the reliability and skill of the probabilistic forecasts (Hamill et

119   al. 2004, Hamill et al. 2006, Hamill and Whitaker 2007, Wilks and Hamill 2007,

120   Hamill et al. 2008).   A drawback of using reforecasts is that a forecast time series

121    spanning many years or even decades may be necessary to produce a sufficiently

122    large sample to adjust for systematic errors in rare-event forecasts.  Since forecast

123    models are frequently updated, which may change the systematic error

124    characteristics, either a forecast model must be frozen once a reforecast data set has

125    been generated, or a new reforecast data set must be generated every time the

126    modeling system changes significantly.  Hence, reforecasting can be computationally

127    expensive and can restrict the ability of a forecast center to upgrade its system

128    rapidly.   Recently, statistical post-processing methods have been the subject of

129    much investigation (Gneiting et al. 2005, Raftery et al. 2005, Sloughter et al. 2007,

130    Wilson et al. 2007, Vannitsem and Nicolis 2008, Glahn et al. 2009, Bao et al. 2010).

131         To date, however, there have been no systematic comparisons of multi-

132    model and reforecast-calibrated PQPFs verified over a large enough area and a long

133    enough period of time to confidently assess the relative strengths and weaknesses

134    of these two approaches.  This study attempts to provide such a comparison for this

135    high-impact forecast parameter.  Using TIGGE forecast data from the US National

136    Centers for Environmental Prediction (NCEP), the Canadian Meteorological Centre

137    (CMC), the United Kingdom Met Office (UKMO), and the European Centre for

138    Medium-Range Weather Forecasts (ECMWF), multi-model ensemble 24-h

139    accumulated probabilistic forecasts of precipitation were generated and then

140    compared against ECMWF forecasts that were statistically adjusted using their

141    reforecast data set.  The comparison was performed over the contiguous US

142    (CONUS) during the period July-October 2010.   Statistical adjustments were also

143     attempted with multi-model forecasts, trained on the previous 30 days of forecasts

144     and analyses.

145         Below, section 2 describes the data sets used in this experiment, the

146     verification methodology, and the statistical post-processing method. Section 3

147     provides results, and section 4 some conclusions.

148
149     2. **Data sets and methods**.

150
151     *a.  Analysis data used.*

152         A recently created precipitation data set, NCEP's Climatology-Calibrated

153     Precipitation Analysis (CCPA), was used for verification.   The CCPA attempts to

154     combine the relative advantages of the 4-km, hourly NCEP Stage-IV precipitation

155     analysis (Lin and Mitchell 2005), and the daily, 0.25-degree NCEP Climate Prediction

156     Center (CPC) Unified Precipitation Analysis (Higgins et al. 1996).  The former is

157     based on gauge and radar data, the latter solely on gauge data.  A disadvantage of

158     the Stage-IV product is that it may inherit some of the biases due to the estimation

159     of rainfall from radars.  A disadvantage of the CPC product is that there are areas of

160     the US that are only sparsely covered by gauge data.   The CCPA analysis regressed

161     the Stage-IV analysis (the predictor) to the CPC analysis (the predictand), thereby

162     reducing bias with respect to the in-situ observations.  In several of the driest

163     locations in the western US, the CCPA analysis was set to missing, for the regression

164     analysis was untrustworthy and singular due to no precipitation in either analysis

165     product.  In such cases, the CCPA analysis for this study was simply replaced with

166     the Stage-IV analysis.   For our purposes, we used CCPA analyses that also were

7

167  upscaled to 1 degree and accumulated over a 24-h period in a manner that

168  preserved total precipitation, similar to the "remapping" procedure described in

169  Accadia et al. (2003).  The CCPA analyses were available from 2002 – current, a

170  shorter period than the ECMWF reforecasts, thus limiting the amount of training

171  data that could be used in the statistical post-processing.

172
173  *b.  Forecast and reforecast model data.*

174       For this experiment, 20 perturbed member forecasts of 24-h accumulated

175  precipitation were extracted from the UKMO, CMC, NCEP, and ECMWF ensemble

176  systems archived in the TIGGE database at ECMWF.  Probabilities were calculated

177  directly from the ensemble relative frequency, referred to as "raw" probabilities

178  henceforth.  The forecast period was July to October 2010; only 00 UTC initial time

179  forecasts were extracted in order to allow comparison with a post-processed

180  forecasts using ECMWF's reforecasts, which were generated only from 00 UTC

181  initial conditions.  Daily forecasts of 24-h accumulated precipitation were examined

182  from +1 to +5 day lead.  Regardless of the original model resolution, all centers'

183  forecasts were bi-linearly interpolated to a 1-degree latitude-longitude grid

184  covering the CONUS using ECMWF's TIGGE portal software.   ECMWF's

185  interpolation procedure  set the amount to zero if there was no precipitation at the

186  nearest neighboring point and the interpolated value was less than 0.05 mm.  No

187  control forecasts were included, just the forecasts from the perturbed initial

188  conditions.  Other forecast centers' contributions to the TIGGE archive were not

189  used here for various reasons, such as the unavailability of 00 UTC ensemble

190  forecasts from the Japan Meteorological Agency.  For size consistency and to

191    facilitate skill comparisons, only the first 20 of the full 50 ECMWF member forecasts

192    were used in the generation of the multi-model ensemble, though the 50-member

193    ECMWF forecasts were evaluated for skill and reliability.   More detailed

194    descriptions of the configuration of these four ensemble systems are described in

195    Appendix 1.

196          When calibrating ECMWF data with reforecasts, the 5-member weekly

197    reforecasts precipitation data were extracted from ECMWF's weekly reforecast

198    archive (Hagedorn 2008) and similarly interpolated to the 1-degree grid.   The

199    control reforecasts were initialized from the ERA-Interim reanalysis (Dee et al.

200    2011), which used version Cy31r2 of the ECMWF Integrated Forecast System (IFS)

201    in the data assimilation process. The 2010 real-time ensemble forecasts and the

202    reforecasts were then run using IFS model version Cy36r2 (more detail is provided

203    in appendix A).

204          The four perturbed initial conditions for the reforecasts were generated with

205    a combination of their singular-vector approach (Molteni et al. 1996, Barkmeijer et

206    al. 1999) and their "ensembles of data assimilations" or "EDA" (Isaksen 2010) that

207    used a cycled, reduced resolution 4D-Var and perturbed observations.  However, for

208    initialization of the reforecasts, ECMWF used the EDA perturbations from 2010 and

209    applied them to the 2002-2009 data rather than running the EDA during the

210    reforecast period.  To apply EDA to dates in the past would have been

211    computationally expensive, but having not done so may have resulted in the

212    perturbations in the reforecast having less flow-dependent character, possibly

213    making them somewhat statistically inconsistent with ECMWF's real-time ensemble

214    forecasts.

215         Since precipitation analysis data was only available for the period from 2002

216    forward, the training data was limited to the reforecasts for period of June to

217    November, 2002-2009, or 8 years.  To limit the possible deleterious effects of

218    seasonal biases in the forecast model, only the reforecast data for the month of

219    interest and the month before and after were used.  For example, when calibrating

220    July forecasts, June-July-August reforecast data was used.   With reforecasts

221    generated once per week, this typically meant there were ~13 once-weekly samples

222    $\times$ 8 years = 104 samples.   Tthis was in many cases an insufficient sample size, so

223    data from other grid points were used to increase  the training sample size

224    (appendix B).

225
226    *c. Statistical post-processing methodology.*

227         The extended logistic regression (ELR) approach of Wilks (2009) was used

228    here, a procedure that permitted the development of a single regression equation

229    that was suitable for predicting probabilities of exceeding any precipitation amount.

230    The probability was estimated with a function of the form

231         $$p = \frac{\exp[f(\mathbf{x})]}{1+\exp[f(\mathbf{x})]},$$    (1)

232    where $f(\mathbf{x})$ was a linear function of the predictor variables.  In this case, the

233    predictors were (a) the ensemble-mean forecast $\bar{x}$ raised to the 0.4 power, (b) the

234    product of (a) and the variance $\sigma^2$ to the 0.4 power, and (c) the precipitation event

235    threshold $T$ raised to the 0.4 power.  The linear function was thus

236 $$f(\mathbf{x}) = b_0 + b_1 \bar{x}^{0.4} + b_2 \bar{x}^{0.4} \sigma^{2 \times 0.4} + b_3 T^{0.4}. \qquad (2)$$

237 The choice of these predictors was arrived at through some trial and error. The

238 power transformation of the predictors helped make the input data somewhat more

239 normally distributed. The probabilistic forecast skill was also only mildly

240 dependent on the inclusion/exclusion of the predictor with the product of the

241 transformed mean and variance. Skill was also only slightly dependent on the

242 power of the transform, with 0.4 providing an approximate minimum. Previous

243 values of power transformations in the literature have ranged from ½ in Hamill and

244 Whitaker (2006) and Schmeits and Kok (2010), 1/3 in Sloughter et al. (2007), and

245 ¼ in Hamill et al. (2008) and Roulin and Vannitsem (2011). The use of the product

246 of the ensemble mean and variance follows Wilks and Hamill (2007). The additional

247 predictor incorporating $T$ permitted the single regression equation to be used to

248 predict probabilities across the range of possible amounts. A disadvantage of this

249 ELR approach (as opposed to approaches such as the analog approach discussed in

250 Hamill and Whitaker (2006)) was that this algorithm was not able to correct for

251 possible position biases in forecast features.

252      ELR was applied both to calibrating real-time multi-model forecasts and to

253 calibrating ECMWF forecasts alone using the weekly reforecasts. It was found that

254 forecast skill increased if some method was applied to increase the modest training

255 sample sizes. A discussion of how sample sizes were augmented using data from

256 other nearby forecast grid points is provided in Appendix 2.

257      Roulin and Vannitsem (2011) noted that since the ECMWF reforecast size (5

258 members) was smaller than the operational ensemble size (50 members; or in the

259    case here, 20 members selected from the 50), the regression coefficients may be

260    somewhat biased when trained with a smaller ensemble compared to what they

261    would be were they trained with a larger ensemble.  Hence, when the coefficients

262    are used to correct the larger real-time ensemble, they may produce somewhat

263    biased probabilistic forecasts.   They adjusted the values of the 5-member ensemble

264    training data to better estimate the values that would be obtained with the larger

265    real-time ensemble.  An analogous approach was tried here but did not improve the

266    forecast skill.  The results discussed below will omit this adjustment.

267
268    *d. Verification methods.*

269         The primary verification methods used here were Brier Skill Scores (*BSS*),

270    continuous ranked probability skill scores (*CRPSS*), and reliability diagrams (Wilks

271    2006).   The *BSS* and *CRPSS* as conventionally calculated (see section 7.4.2 of Wilks

272    (2006)) can exaggerate forecast skill, attributing skill to variations in climatological

273    event probabilities. Thus, the procedures suggested in Hamill and Juras (2006) were

274    used here to avoid this.

275         To calculate the *BSS*, the score was calculated separately for subsets of points

276    that had more uniform climatological probabilities.  The overall *BSS* was the average

277    of the skill scores over these subsets.  The specific procedure was as follows.  Using

278    the 1-degree precipitation analysis data from 2002-2009, for each month the

279    climatological probability of a given precipitation event was estimated from the

280    observed frequency.   For a given event such as > 1 mm (24 h)$^{-1}$, the $n_s$ grid points

281    within the CONUS were sorted from lowest to highest event probability. The sorted

282    points were then divided into $k$=6 classes, with the lowest bin containing the $\sim n_s/6$

283     grid points with the lowest event probabilities, the highest bin containing the $n_s/6$

284     points with the highest probabilities, and so on (Fig. 1). Let $\mathbf{BS}^{f1} = \left[ \mathbf{bs}_1^{f1}, \ldots, \mathbf{bs}_6^{f1} \right]$

285     denote a matrix of Brier scores for forecast model *f1*, where $\mathbf{bs}_i^{f1}$ was a $n_d -$

286     dimensional (= 123, the number of case days here) column vector of average Brier

287     scores for the points in the $i^{\text{th}}$ class and for forecast model *f1*.  An element of this

288     vector thus provided the average Brier Score for all of the grid points in the $i^{\text{th}}$ class

289     on a particular day; the samples were weighted by the cosine of their latitude to

290     account for differences in grid box size.  The average over the 123 case days

291     produced a vector $\overline{\mathbf{bs}}^{f1} = \left[ \overline{bs}_1^{f1}, \ldots, \overline{bs}_6^{f1} \right]$.  Similarly, for climatology there was an

292     array of Brier scores, $\mathbf{BS}^c = \left[ \mathbf{bs}_1^c, \ldots, \mathbf{bs}_6^c \right]$ and a vector of their averages over the

293     123 days, $\overline{\mathbf{bs}}^c = \left[ \overline{bs}_1^c, \ldots, \overline{bs}_6^c \right]$.  Following Hamill and Juras (2006) eq. (9), the

294     overall *BSS* for model *f1* was then calculated as

295     $$BSS = \sum_{k=1}^{6} \frac{1}{6} \left( 1 - \frac{\overline{bs}_k^{f1}}{\overline{bs}_k^c} \right). \tag{3}$$

296     The boundaries between the classes were calculated independently for each event,

297     so it was possible that a given grid point may have been assigned to different classes

298     when evaluating, say, the 1- and 10-mm *BSS*s.

299          *BSS* confidence intervals were estimated using the paired block bootstrap

300     approach of Hamill (1999).  The input data to the bootstrap approach consisted of

301     arrays of $\mathbf{BS}^{f1}$ and $\mathbf{BS}^{f2}$ for two competing models, *f1* and *f2*, as well as $\mathbf{BS}^c$.  Let

302     $\mathbf{bs}^{f1}(d) = \left[ bs_1^{f1}(d), \ldots, bs_6^{f1}(d) \right]$ be the vector of forecast scores on the $d^{\text{th}}$ case day,

303  and similarly $\mathbf{bs}^{f2}(d)$ the vector for forecast model *f2*. The daily differences in Brier

304  scores, $\mathbf{bs}^{f1}(d) - \mathbf{bs}^{f2}(d)$ were determined to be approximately statistically

305  independent of $\mathbf{bs}^{f1}(d+1) - \mathbf{bs}^{f2}(d+1)$, with 1-day lagged rank correlations ranging

306  from 0.08 for 1-day forecasts to 0.21 for 5-day forecasts.  Thus, the data was judged

307  to be amenable to a paired resampling strategy, with these distinct vector blocks of

308  data for each day.   The following process was then repeated 10,000 times.  For each

309  of the 123 days, a random uniform number between 0 and 1 was generated.  If the

310  number was greater than 0.5, $\mathbf{bs}^{f1}(d)$ was randomly selected for inclusion in sample

311  1, $\mathbf{bs}^{f2}(d)$ was selected for inclusion in sample 2, and vice versa if the number was

312  less than or equal to 0.5.   The vector of average Brier scores for samples *s1* and *s2*

313  were then calculated, $\overline{\mathbf{bs}}^{s1}$ and $\overline{\mathbf{bs}}^{s2}$.  The *BSS* for samples 1 and 2 were generated

314  via eq. (1), and the difference between the *BSS*s for the two samples was noted.  The

315  confidence intervals are the 5th and 95th percentiles of the difference between the

316  *BSS*s of the two samples from the distribution generated through the 10,000

317  iterations.

318       These block bootstrap confidence intervals should be regarded as

319  approximations.  An assumption underlying this process is that there were 123

320  independent data samples.  However, $\mathbf{bs}^{f1}(d)$ and $\mathbf{bs}^{f1}(d+1)$ were slightly

321  correlated as discussed above, especially for the longer-lead forecasts, which will

322  contribute a slight over-estimate of the effective sample size and thus underestimate

323  of the confidence interval.   On the other hand, data from grid points across the

324  CONUS were aggregated in this procedure and thereafter considered as a single

325    block.  In reality there may be far more than one independent sample spanning the

326    CONUS, thus leading to an under-estimate of sample size and consequent

327    overestimate of the confidence interval in this approach.  Note also that for

328    simplicity of presentation, the skill diagrams will show only one set of confidence

329    intervals, e.g., between NCEP and ECMWF forecasts.  Slightly smaller confidence

330    intervals could be expected were they computed using ECMWF and CMC forecasts,

331    given their more similar skills.

332         In order to make sure that the *CRPSS* did not excessively reflect the skill of

333    the climatologically wet grid points, an alternative to the standard method of *CRPSS*

334    calculation was followed here.  This method is analogous to how *CRPSS* would be

335    computed if the forecasts were of probabilities of exceedance of various quantiles.

336    An example of such a forecast product expressed in quantiles are NCEP/CPC's 6-10

337    day and 8-14 day forecasts (e.g.,

338    http://www.cpc.ncep.noaa.gov/products/predictions/610day/), which provide

339    probabilities of below-normal/near-normal/above-normal temperature and

340    precipitation, i.e., probabilities for the $< 1/3$ and $\geq 2/3$ quantiles.  In this alternative

341    method of calculation, the *CRPS* at a given grid point was *not* computed by

342    integrating differences between observed and forecast cumulative distribution

343    functions (CDFs) *over a range of precipitation values* (the standard method).  Instead,

344    the differences between observed and forecast CDFs were integrated *over the*

345    *percentiles of the CDF*, which were determined separately for each model grid point

346    and each month.  Specifically, given $n_d$ case days, for the $s = 1, \ldots, n_d \times n_s$ samples, let

347    $\mathbf{q}^s = \left[ q_1^s, \ldots, q_{20}^s \right]$ be the 20-dimensional vector of the precipitation quantiles

15

348    associated with the 2.5[th], 7.5[th], ..., 97.5[th] percentiles of the climatological *CDF* for that

349    point and that month.  The average forecast *CRPS_f* was determined by integrating in

350    steps of 5 percent:

$$CRPS_f = \frac{\sum_{s=1}^{n_d \times n_s} \cos(\phi_s) \sum_{iq=1}^{20} 0.05 \times \left[ F^s\left(q_{iq}^s\right) - O^s\left(q_{iq}^s\right) \right]^2}{\sum_{s=1}^{n_d \times n_s} \cos(\phi_s)} \tag{4}$$

,

352    where $F^s\left(q_{iq}^s\right)$ represents the forecast's *CDF* for the $s$[th] sample evaluated at the $q_{iq}^s$

353    quantile, and $O^s\left(q_{iq}^s\right)$ represents the same, but for the observed (analyzed). $\phi_s$ is the

354    latitude of the grid box, the cosine factor accounting for latitudinal variations in grid

355    box size. For raw ensemble guidance, $F^s\left(q_{iq}^s\right)$ was directly computed from the

356    ensemble relative frequency.  For example, if 5 of 20 members had values exceeding

357    $q_{iq}^s$, then $F^s\left(q_{iq}^s\right) = 0.75$.  For post-processed forecasts, $F^s\left(q_{iq}^s\right)$ was determined by

358    numerical integration of eqs. (1) and (2).  For the observed CDF, the analyzed state

359    was assumed perfect, i.e., no analysis errors were incorporated, so the analyzed CDF

360    was a Heaviside function, 0 at the quantiles less than the analyzed value, 1 at

361    quantiles greater than or equal to the analyzed value.  The *CRPS* of the climatological

362    forecast, *CRPS_c* , was calculated as in eq. (4), but substituting the climatological CDF

363    for the forecast CDF.  Finally, the overall skill score was calculated as *CRPSS* = 1. -

364    *CRPS_f* / *CRPS_c*.  As with the *BSS*, a paired block bootstrap approach was used to

365    estimate the confidence intervals.

366         Two other common verification statistics were also used, root-mean-square

367    (RMS) errors, and bias, the average forecast divided by the average analyzed

368    amount.

369

370    3. **Results**.

371    *a. Properties of forecasts from the individual centers.*

372         Before considering the multi-model and ECMWF reforecast-calibrated

373    forecast properties, let us consider the verification of PQPFs from the individual

374    centers.  Figure 2 shows > 1 mm $(24 \text{ h})^{-1}$ and > 10-mm $(24 \text{ h})^{-1}$ event *BSS* and *CRPSS*.

375    ECMWF generally produced the most skillful raw precipitation PQPFs. Depending on

376    the metric, either NCEP or UKMO produced the least skillful forecasts.

377         Interestingly, though UKMO forecasts appeared to be generally more skillful

378    than NCEP forecasts in *BSS*, they appeared to be consistently worse in *CRPSS*.  This

379    was a consequence of the *CRPSS* verification algorithm as implemented here, which

380    attempted to equally weight the *CRPSS* at all grid points, irrespective of whether the

381    climatological event probability was extremely high or extremely low.  The

382    conventionally calculated *CRPSS* is dominated by the performance of the forecasts in

383    the climatologically wet areas (Hamill and Juras 2006).  There is inherently greater

384    climatological variance of precipitation for the wet regions, and associated with this

385    there are generally much larger *CRPS* values than in dry regions.  Consequently,

386    when evaluated over many grid points, the conventionally calculated *CRPS* and

387    hence the *CRPSS* are dominated by the performance at the wetter points.  Figure 3

388    shows maps of the day +3 *CRPSS* scores (see the online appendix for *CRPSS* maps for

389    the other lead times).  The UKMO forecasts had negative skill in the extremely dry

390    regions of the western US.  The RMS errors of the ensemble-mean forecasts in the

391    dry regions of all the models were very small and relatively similar (Fig. 4a; for

392    other lead times, see the online appendix).  However, the UKMO forecasts exhibited

17

393 a large moist bias in the climatologically dry regions (Fig. 4b), which resulted in a

394 very large over-forecast of probabilities and poor skill for those points.   This was

395 apparently due to a drizzle over-forecast bias in that version of the UKMO's forecast

396 model (D. Barker, personal communication, 2011).

397   Figure 4b also illustrates some other interesting characteristics of the

398 ensemble systems.   NCEP over-forecasted rainfall for the grid points and dates

399 where the climatological probability was already quite high.  CMC forecasts were

400 also biased, exhibiting a moist bias at the lowest climatological probabilities but dry

401 biases for most of the rest of the larger climatological probabilities.  ECMWF

402 forecasts were the least biased, with a moderate over-forecast bias at the low

403 climatological probabilities.

404   Figure 5 provides reliability diagrams of day +3 forecasts of the > 10-mm (24

405 h)$^{-1}$ event.  Other reliability diagrams for other lead times and for the > 1-mm (24 h)$^{-}$

406 $^1$ event are available in the online appendix.  CMC forecasts were generally the most

407 reliable, though they were not as sharp as the ECMWF forecasts and hence had a

408 lower *BSS*. UKMO and NCEP forecasts were much less reliable, though NCEP

409 forecasts were slightly sharper than the others.  ECMWF 50-member forecasts were

410 slightly more reliable and skillful than their 20-member subset.

411   In subjective analyses of individual forecasts, it appeared that several of the

412 forecast models had subtle systematic northward biases in the northern central US.

413 Figure 6 shows the 10-mm observed contour and the 0.5 probability contour for the

414 > 10-mm (24 h)$^{-1}$ event from the day +3 ECMWF forecasts.  Here, the 25 cases with

415 the largest areal coverage of observed precipitation between 105° and 80° west

416    longitude and 35° and 50° north latitude were chosen.    Similar plots for the other

417    forecast models are included in the online appendix.

418
419    *b. Properties of multi-model and statistically post-processed forecasts.*

420            Before considering verification scores, consider first two actual forecast

421    cases, presented in Figs. 6 and 7, showing probabilities from the 20-member

422    ensembles and from the 80-member multi-model ensemble.  The first case, covering

423    the 24-h period ending 00 UTC 21 July 2010, illustrates that sometimes the forecast

424    models could be overly similar to each other.  Here all the forecast precipitation

425    areas were significantly north of the observed area.  A multi-model forecast would

426    not be expected to provide much benefit in such a situation.  Figure 7 shows the

427    same type of plot, but for 24-h period ending 00 UTC 7 August 2010.  Here the multi-

428    model forecast provided some improvement.   On this day the CMC and UKMO areas

429    of high probabilities were too far north, the NCEP area too far south, but the higher

430    probabilities in the multi-model forecasts were more coincident with the analyzed

431    regions exceeding 10 mm.  Most of the area with greater than 10 mm in the analysis

432    were covered by nonzero multi-model probabilities.  More generally, when there

433    was some diversity of positions in the multi-model forecasts, this often allowed the

434    forecast to avoid being inappropriately sharp.

435            Figure 9 provides *BSS*s and *CRPSS* for the multi-model and the post-

436    processed forecasts.  For the light precipitation forecasts (> 1.0 mm (24 h)$^{-1}$, Fig. 9a),

437    the multi-model forecasts improved the skill by approximately +1 day relative to

438    ECMWF at the earliest lead times; a +2 day multi-model forecast could now be made

439    as skillfully as a +1 day ECMWF forecast.  The improvement in skill was a more

440    modest ~ +0.3 days at the longer forecast lead times.  The calibrated multi-model

441    forecast product improved skill over the basic multi-model forecast by a tiny

442    amount at day +1 but degraded the skill after day +3.  This is consistent with

443    previous results; at the longer lead times, the growth of errors makes it more

444    difficult to differentiate the model bias from the chaotically induced errors with

445    short training data sets (Hamill et al. 2004).   The improvement from reforecast-

446    based post-processing over the raw ECMWF system was much smaller than the

447    improvement from single to multi-model and was even slightly negative at the day

448    +5 lead.   Reasons for the less impressive performance of reforecast calibration than

449    in previous studies will be discussed at the end of this section.

450        More impressive increases in skill were evident for the > 10-mm $(24\ \text{h})^{-1}$

451    event.  Both the reforecast-based calibration and the multi-model approach

452    increased forecast skill by an equivalent of up to +2 days of additional lead time.

453    Again, the calibration of the multi-model forecasts provided modest improvement at

454    the early leads and degradation at the longer leads relative to the unprocessed

455    multi-model.

456        Measured in *CRPSS*, the multi-model forecasts produced the most skillful

457    forecasts, exceeding the skill of reforecast-calibrated ECMWF forecasts by a small

458    amount.  Consider now *where* the forecasts were improved or degraded by the

459    various approaches.  Figure 10 provides maps of the day +3 *CRPSS*; maps for other

460    lead times are in the online appendix.  The patterns of multi-model skill are rather

461    similar to those of the most skillful ensemble system, ECMWF (Fig. 3a).  The

20

462    reforecast-calibrated ECMWF forecasts appear to have increased the skill most

463    notably in the driest regions of the western US.

464        Figure 11 shows day +3 > 10-mm $(24 \text{ h})^{-1}$ event reliability diagrams for the

465    multi-model, the calibrated multi-model, and reforecast-calibrated ECMWF PQPFs.

466    The raw multi-model PQPFs were slightly more reliable than any of the PQPFs from

467    the individual centers (Fig. 5) and retained a slight over-forecast bias at the higher

468    probabilities.  The improvements in reliability were more substantial than for the >

469    1-mm $(24 \text{ h})^{-1}$ event; see diagrams in the online appendix.  The reforecast-calibrated

470    PQPFs exhibited a slight under-forecast bias and were not as sharp as those from

471    the multi-model forecasts.  Was this due to some inhomogeneity between the 2002-

472    2009 training data and the 2010 real-time forecasts?  Figure 12 shows that there

473    were fewer large forecast busts in 2010 than there were in 2002 or 2006.  When the

474    regression analysis from 2002-2009 data was applied to correct the 2010 forecasts,

475    the assumption was that the 2010 forecasts would be equally unskillful.  In fact they

476    were better, and as a consequence the post-processed forecasts were less sharp

477    than they could have been.  Though it was not attempted here, it might be possible

478    to apply ad-hoc corrections to the training data to improve the regression analysis.

479    Perhaps a slight adjustment of the training data ensemble mean toward the

480    analyzed data would make its accuracy more closely resemble that of the 2010 data,

481    sharpening and making the ELR forecasts more reliable and skillful.

482        Figure 13 shows the multi-model areal coverage of the 0.5 probability

483    contours for the > 10 mm $(24 \text{ h})^{-1}$ event for selected cases; these should be

484    compared with Fig. 6 for ECMWF-only PQPFs.  Figure 14 also shows the areal

485    coverage, but for reforecast-calibrated ECMWF PQPFs.  The areal coverage was only

486    slightly smaller for the multi-model PQPFs than it was for the ECMWF PQPFs,

487    illustrating that the multi-model forecasts did not lose a tremendous amount of

488    sharpness (coverage of greater than 0.5 probability being a proxy for sharpness

489    here).  In comparison, the reforecast-calibrated PQPFs in Fig. 14 show a marked

490    decrease in the areal coverage; many grid points with probability $p > 0.5$ in the raw

491    ECMWF PQPF had $p < 0.5$ after calibration.  Figures 15 and 16 show for the cases

492    plotted in Figs. 7 and 8 a bit more detail on what happened with typical multi-model

493    and reforecast-calibrated PQPFs.  The multi-model forecasts retained their

494    sharpness, but not always desirably so.  For example, in Fig. 15, the multi-model

495    forecasts retain relatively high probabilities in eastern Iowa and northern Illinois,

496    whereas the analyzed area was displaced further south.  The reforecast-calibrated

497    PQPFs decreased the areal coverage of high probabilities, appropriately so in this

498    case, reducing the false alarms.  However, as seen in inspection of Figs. 13-14, there

499    were many cases when the sharpness retained in the multi-model forecasts was

500    desirable.

501        The results exhibited here with reforecast calibration were not as impressive

502    as they have been in previous studies, e.g., Hamill and Whitaker (2006) and Hamill

503    et al. (2008).  There are at least four reasons for this.  First, the training data was not

504    as accurate as the real-time data in this application (Fig. 12), and this inhomogeneity

505    degraded the regression analysis.  This may have been due to less accurate initial

506    conditions (ERA-Interim for the reforecast, operational 4D-Var for the real-time

507    forecasts) and because the reforecast ensemble was initialized with perturbations

22

508    that were constructed with approximations different from those in the real-time

509    forecasts (section 2b).  The second reason is that gratifying improvements have

510    been made to models and EPSs so that they produce more skillful and reliable

511    forecasts than they did even in the recent past; it's tougher to improve upon

512    ECMWF's 2010's model output than its 2005 model output.  The third reason is that

513    even with the use of ECMWF's reforecasts, there really was a limited training data

514    set in this study, here due to the unavailability of precipitation analyses prior to

515    2002 and the unavailability of reforecast data more frequently than once per week.

516    The fourth reason is that in prior studies, the ensemble forecasts (at coarse

517    resolution) were evaluated against analysis data at finer resolution, so that the

518    reforecast calibration process was also producing a statistical downscaling.  This

519    point is worth keeping in mind when considering the relative merits of reforecast

520    calibration vs. multi-model approaches.  If the desired output is forecast data at the

521    grid scale, multi-models may have substantial appeal.  If the desired output is point

522    data or high-resolution gridded data, the statistical downscaling is more

523    straightforward when reforecasts are used.

524          Overall, the impressive skill improvements provide evidence for the merit of

525    both multi-model ensemble and reforecast approaches.   Should other forecast

526    centers share precipitation ensemble data, large gains in probabilistic precipitation

527    forecast skill are possible for little more than the cost of data transmission and

528    storage.  Alternatively, should any one center produce and utilize reforecasts, they

529    can improve their own forecasts significantly, assuming a comparably long time

530    series of observations or analyses are available.  The improvement here noted with

531    reforecasts may have also been modest because the training data was limited on

532    account of a short time series of analyses, dating back to only 2002; only around

533    40% of the available reforecast data was used.

534
535    **4. Conclusions**.

536         This article examined probabilistic multi-model weather forecasts of

537    precipitation over the CONUS and the relative advantages and disadvantages of

538    these forecasts when compared to statistically post-processed ECMWF forecasts.

539    20-member forecasts were extracted from the ECMWF, NCEP, UKMO, and CMC

540    global ensemble systems at 1-degree resolution between June and October 2010.

541    Daily 24-h accumulated probabilistic precipitation forecasts were generated from

542    the subsequent 80-member ensemble for lead times of +1 to +5 days and compared

543    to gridded precipitation analyses.  Two statistically post-processed products were

544    also evaluated, the first being multi-model forecasts that were adjusted using

545    extended logistic regression and that were trained on the previous 30 days of

546    forecasts and analyses.  The second was ECMWF forecasts, which were statistically

547    adjusted using forecast/analysis data for the period 2002-2009, the time period

548    when both reforecasts and analyses were available.

549         Considering first the skill of forecasts from the individual EPSs, ECMWF

550    forecasts generally were the most skillful in terms of Brier skill scores and the

551    continuous ranked probability skill score.  CMC forecasts were the most reliable but

552    the least sharp, while NCEP and UKMO forecasts were more sharp but less reliable.

553         Multi-model probabilistic forecast products were substantially more skillful

554    than the best of the individual centers' probabilistic forecasts.   The improvement

555    was approximately an extra +0.5 to +1 day of forecast lead time for light

556    precipitation events and as much as +2 days for heavier precipitation events.  The

557    reforecast-calibrated ECMWF forecasts exhibited more skill and reliability

558    improvement at the > 10-mm $(24\ h)^{-1}$ event as they did at the > 1-mm $(24\ h)^{-1}$ event.

559    Relative to the multi-model forecasts, the reforecast-calibrated skills were similar

560    for the > 10 mm $(24\ h)^{-1}$ event, but the reforecast-calibrated was more reliable while

561    the multi-model was sharper.

562        The results exhibited here with reforecast calibration were not as impressive

563    as they have been in previous studies.  There were at least four reasons for the

564    lessened improvement of reforecast calibration here.  First, the reforecast training

565    data was shown to be not as accurate as the real-time data in this application.

566    Second, gratifying improvements have been made to models and EPSs in the last few

567    years; it's tougher to improve upon ECMWF's 2010's model output than its 2005

568    model output.  Third, limited training data set was available for this study.  Fourth,

569    prior studies were performed at higher resolution and produced a statistical

570    downscaling that the coarser raw forecasts could not accomplish.

571        I was pleasantly surprised by the magnitude of skill improvements

572    demonstrated here from multi-model ensembles, improvements which were larger

573    than those seen with 2-meter temperatures (Hagedorn et al. 2011).  From our own

574    experience, however, I recommend some caution against broadly generalizing these

575    results to any multi-model ensemble system.  This study examined a combination of

576    data from four mature EPSs based on mature models and assimilation systems.

577    Each center's system has been refined through the collective efforts of hundreds if

578    not thousands of person-years of research and development.   A combination of less

579    developed EPSs may not provide nearly the same gratifying result.

580         Nonetheless, these results demonstrate the potential value of multi-model

581    ensembles.   The THORPEX program, organized by the World Meteorological

582    Organization, has promoted the concept of a multi-model based "Global Interactive

583    Forecast System" (Bougeault et al. 2010),  whereby the operational centers share

584    data that will facilitate the production of multi-model products for high-impact

585    weather events.  This study provides additional evidence for the validity and the

586    potential benefits of such a system.  Currently several centers have restrictive data

587    policies; full access to their data is reserved for paying customers, and those

588    customers cannot thereafter share the data they purchased.  Perhaps the approach

589    embraced in the US and Canada will be followed by other centers worldwide, for the

590    mutual benefit of all.  In the US and Canada, the data is effectively free since the

591    research, development, and production were funded by public taxpayer funds.

592         Finally, can we all have "the best of both worlds?"  That is, will NWP centers

593    both agree to share their ensemble data freely and internationally in real time, and

594    will they produce reforecast data sets so that each model can be calibrated to

595    remove systematic errors prior to their combination?  There is evidence that such

596    approaches will provide substantial benefit.  The climate community is working on

597    sharing multi-model information and hindcasts to facilitate the error correction for

598    intra-seasonal and seasonal forecasts.  For weather and weather-to-climate

599    applications, there have also been successful demonstrations of multi-model

600    calibrated forecasts (Vislocky and Fritsch 1995, Whitaker et al. 2006).   NOAA is

601 currently developing a new reforecast data set for its global ensemble prediction

602 system, and I hope that other centers will be inspired to do so as well.

603
604
605 **Acknowledgments**

614

615 **Appendix 1.**

616 Here are additional details on the forecast models and ensemble systems used in

617 this experiment.

618
619 *a. NCEP*

620     NCEP used the Global Forecast System (GFS) model in their ensemble system

621 at T190L28 resolution. A lengthier description of the physical packages used in this

622 model were described in Hamill et al. (2011). A description of the GFS model is

623 available from the NCEP Environmental Modeling Center (EMC), with changes as of

624 2003 described at www.emc.ncep.noaa.gov/gmb/moorthi/gam.html.

625     The control initial condition around which the perturbed initial conditions

626     were centered was produced by the T382 Global Statistical Interpolation (GSI)

627     analysis (Kleist et al. 2009) at T384L64 resolution.  Perturbed initial conditions

628     were generated with the ensemble transform with rescaling technique of Wei et al.

629     (2008). Stochastic perturbations were included, following  Hou et al. (2008). More

630     details on changes to the NCEP ensemble system can be found at

631     http://www.emc.ncep.noaa.gov/gmb/yzhu/html/ENS_IMP.html.

632
633     *b.  Canadian Meteorological Centre*

634     The CMC EPS used the Global Environmental Multiscale Model, a primitive

635     equation model with a terrain-following pressure vertical coordinate.  Further

636     documentation on the GEM model can be found at

637     http://collaboration.cmc.ec.gc.ca/science/rpn/gef_html_public/DOCUMENTATION/

638     GENERAL/general.html and in Charron et al. (2010). The CMC ensemble system

639     used a horizontal computational grid of 400x200 grid points, or approximately 0.9

640     degrees, and 28 vertical levels.  The EnKF initial conditions were used, following

641     Charron et al. (2010) and Houtekamer et al. (2009) and references therein.  The 20

642     forecast ensemble members used a variety of perturbed physics; changing gravity

643     wave drag parameters, land-surface process type, condensation scheme type,

644     convection scheme type, shallow convection scheme type, mixing-length

645     formulation, and turbulent vertical diffusion parameter.  More details on these are

646     provided at http://www.weatheroffice.gc.ca/ensemble/verifs/model_e.html.

647
648     *c. European Centre for Medium-Range Weather Forecasts.*

649
650          The ECMWF EPS used the ECMWF Integrated Forecast System (IFS) model,

651    versions 36r2. Model resolution was T639L62 for both versions; details on the IFS

652    are provided at www.ecmwf.int/research/ifsdocs/.  The changes to the ensemble

653    stochastic treatments in the 8 Sep 2009 implementation are described in Palmer et

654    al. (2009).  The ensemble was initialized with a combination of initial-time and

655    evolved total-energy singular vectors (Buizza and Palmer 1995, Molteni et al. 1996,

656    Barkmeijer et al. 1998, Barkmeijer et al. 1999, Leutbecher 2005) and utilized

657    stochastic perturbations to physical tendencies.  An overview of the ensemble

658    system was provided in Buizza et al. (2007) and references therein.   For

659    consistency with the analysis of other EPSs, only the first 20 perturbed members

660    were used here.

661
662    *e.  United Kingdom Met Office.*

663
664          The UK Met Office (UKMO) ensemble system was "MOGREPS," the Met Office

665    Global and Regional Ensemble Prediction System.  TC track forecasts from this

666    system came from its global component, which was described in Bowler et al. (2008,

667    2009).  The global system was run at a resolution of 0.83° longitude and 0.55°

668    latitude on a regular latitude-longitude grid.  70 vertical levels were employed

669    (Tennant et al. 2011).  Initial condition perturbations were generated from an

670    implementation of the ensemble transform Kalman filter  (Hunt et al. 2006, Bowler

671    et al. 2009).  The mean initial state was generated from the UKMO 4D-Var system

672    (Rawlins et al. 2007). The model included a parameterization of one type of model

673    uncertainty via its stochastic kinetic-energy backscatter scheme, following Shutts

674     (2005) and Tennant et al. (2011).

675

676  **Appendix 2:**

677       This appendix discusses the method used to augment the training sample

678  size used in the regression analyses.  Were only the data at the grid point of interest

679  used for training, when calibrating using the multi-model ensemble using the past

680  30 days of forecasts and analyses, this would provide, of course, only 30 training

681  samples.  Older forecasts could be used, but precipitation biases are often seasonally

682  dependent, so the older data may degrade the results despite augmenting the

683  sample size.  Also, with such a multi-model ensemble, the farther back into the past

684  one seeks training data, the more likely it is that at least one of the models will have

685  had a major upgrade and concomitant change in systematic error characteristics.

686       Despite ECMWF providing a multi-decadal reforecast, in practice the sample

687  sizes were too small here, too. When using the 2002-2009 weekly, 5-member

688  ECMWF reforecasts (including reforecast dates +/- 6 weeks around the week of

689  interest), this provided a total of 13 weeks $\times$ 8 years = 104 samples.  In both cases,

690  these were relatively small samples to estimate four regression parameters, and

691  especially for rare events such as heavy precipitation, experience has shown that

692  larger training samples improved the regression analysis.

693       Hence, following the general philosophy demonstrated and discussed in

694  Hamill et al. (2008) and inspired by the regionalization used in some Model Output

695  Statistics algorithms (Lowry and Glahn 1976), the training data set for a particular

696  grid point was augmented by finding 25 other grid points that had relatively similar

697  climatological analyzed CDFs.  Consider a particular location $(\lambda, \phi)$ at which we seek

698    to augment the sample size, and another location $\left(\lambda^s, \phi^s\right)$ we are considering as a

699    location with suitable supplemental training data. Differences between the analyzed

700    cumulative probabilities at $\left(\lambda, \phi\right)$ and $\left(\lambda^s, \phi^s\right)$ were measured at the 1, 2.5, 5, 10, 25,

701    and 50 mm $(24\ \mathrm{h})^{-1}$ amounts and then weighted by similar respective factors of  [1,

702    2.5, 5, 10, 25, 50].  That is, a cumulative probability difference of 0.1 at 1 mm and

703    0.1/50 at 50 mm were judged to have the same weighted difference (this approach

704    is admittedly somewhat arbitrary, and testing found that the overall calibration

705    results were relatively insensitive to the details of this assumption).  The maximum

706    weighted difference at any of the possible precipitation amounts was then noted for

707    this $\left(\lambda^s, \phi^s\right)$. Having evaluated the maximum of the weighted differences all the grid

708    points less that 8 grid points distant from the grid point of interest $\left(\lambda, \phi\right)$, the 25

709    grid points with the smallest weighted differences were identified, and the training

710    sample for $\left(\lambda, \phi\right)$ was augmented by the forecasts-analysis pairs at these locations.

711    This approach increased sample size, but it's possible that the forecast bias might

712    have been different at the supplemental locations, and hence not an unalloyed

713    benefit.  For more discussion of this, see Hamill et al. (2008, section 3a).

714

715

716
717
718

**719**

**720** **References**
**721**
**722**
**723** Accadia, C., S. Mariani, M. Casaioli, A. Lavagnini, and A. Speranza, 2003: Sensitivity of

**724**         Precipitation Forecast Skill Scores to Bilinear Interpolation and a Simple

**725**         Nearest-Neighbor Average Method on High-Resolution Verification Grids.

**726**         *Weather and Forecasting*, **18,** 918-932.

**727** Bao, L., T. Gneiting, E. P. Grimit, P. Guttorp, and A. E. Raftery, 2010: Bias Correction

**728**         and Bayesian Model Averaging for Ensemble Forecasts of Surface Wind

**729**         Direction. *Monthly Weather Review*, **138,** 1811-1821.

**730** Barkmeijer, J., F. Bouttier, and M. Van Gijzen, 1998: Singular vectors and estimates of

**731**         the analysis-error covariance metric. *Quarterly Journal of the Royal*

**732**         *Meteorological Society*, **124,** 1695-1713.

**733** Barkmeijer, J., R. Buizza, and T. N. Palmer, 1999: 3D-Var Hessian singular vectors

**734**         and their potential use in the ECMWF ensemble prediction system. *Quarterly*

**735**         *Journal of the Royal Meteorological Society*, **125,** 2333-2351.

**736** Berner, J., S.-Y. Ha, J. P. Hacker, A. Fournier, and C. Snyder, 2011: Model uncertainty

**737**         in a mesoscale ensemble prediction system: Stochastic versus multi-physics

**738**         representations. *Monthly Weather Review*, **139,** 1972-1995.

**739** Berner, J., G. J. Shutts, M. Leutbecher, and T. N. Palmer, 2009: A Spectral Stochastic

**740**         Kinetic Energy Backscatter Scheme and Its Impact on Flow-Dependent

**741**         Predictability in the ECMWF Ensemble Prediction System. *Journal of the*

**742**         *Atmospheric Sciences*, **66,** 603-626.

743    Bougeault, P., and Coauthors, 2010: The THORPEX Interactive Grand Global

744        Ensemble. *Bulletin of the American Meteorological Society*, **91,** 1059-1072.

745    Bowler, N. E., A. Arribas, S. E. Beare, K. R. Mylne, and G. J. Shutts, 2009: The local

746        ETKF and SKEB: Upgrades to the MOGREPS short-range ensemble prediction

747        system. *Quarterly Journal of the Royal Meteorological Society*, **135,** 767-776.

748    Bowler, N. E., A. Arribas, K. R. Mylne, K. B. Robertson, and S. E. Beare, 2008: The

749        MOGREPS short-range ensemble prediction system. *Quarterly Journal of the*

750        *Royal Meteorological Society*, **134,** 703-722.

751    Buizza, R., J.-R. Bidlot, N. Wedi, M. Fuentes, M. Hamrud, G. Holt, and F. Vitart, 2007:

752        The new ECMWF VAREPS (Variable Resolution Ensemble Prediction System).

753        *Quarterly Journal of the Royal Meteorological Society*, **133,** 681-695.

754    Buizza, R., M. Miller, and T. N. Palmer, 1999: Stochastic representation of model

755        uncertainties in the ECMWF ensemble prediction system. *Quarterly Journal of*

756        *the Royal Meteorological Society*, **125,** 2887-2908.

757    Buizza, R., and T. N. Palmer, 1995: The Singular-Vector Structure of the Atmospheric

758        Global Circulation. *Journal of the Atmospheric Sciences*, **52,** 1434-1456.

759    Candille, G., 2009: The Multiensemble Approach: The NAEFS Example. *Monthly*

760        *Weather Review*, **137,** 1655-1665.

761    Charron, M., G. r. Pellerin, L. Spacek, P. L. Houtekamer, N. Gagnon, H. L. Mitchell, and

762        L. Michelin, 2010: Toward Random Sampling of Model Error in the Canadian

763        Ensemble Prediction System. *Monthly Weather Review*, **138,** 1877-1901.

764    Dee, D. P., and Coauthors, 2011: The ERA-Interim reanalysis: configuration and

765        performance of the data assimilation system. *Quarterly Journal of the Royal*

766        *Meteorological Society*, **137,** 553-597.

767    Doblas-Reyes, F. J., R. Hagedorn, and T. N. Palmer, 2005: The rationale behind the

768        success of multi-model ensembles in seasonal forecasting – II. Calibration

769        and combination. *Tellus A*, **57,** 234-252.

770    Glahn, B., M. Peroutka, J. Wiedenfeld, J. Wagner, G. Zylstra, B. Schuknecht, and B.

771        Jackson, 2009: MOS Uncertainty Estimates in an Ensemble Framework.

772        *Monthly Weather Review*, **137,** 246-268.

773    Gneiting, T., A. E. Raftery, A. H. Westveld, and T. Goldman, 2005: Calibrated

774        Probabilistic Forecasting Using Ensemble Model Output Statistics and

775        Minimum CRPS Estimation. *Monthly Weather Review*, **133,** 1098-1118.

776    Hagedorn, R., 2008: Using the ECMWF reforecast data set to calibrate EPS

777        reforecasts. *ECMWF Newsletter*, **117,** 8-13.

778    Hagedorn, R., R. Buizza, T. M. Hamill, M. Leutbecher, and T. N. Palmer, 2011:

779        Comparing TIGGE multi-model forecasts with reforecast-calibrated ECMWF

780        ensemble forecasts. *Quarterly Journal of the Royal Meteorological Society*,

781        **accepted pending revision; available from**

782        **martin.leutbecher@ecmwf.int.**

783    Hagedorn, R., F. J. Doblas-Reyes, and T. N. Palmer, 2005: The rationale behind the

784        success of multi-model ensembles in seasonal forecasting – I. Basic concept.

785        *Tellus A*, **57,** 219-233.

786    Hagedorn, R., T. M. Hamill, and J. S. Whitaker, 2008: Probabilistic Forecast

787        Calibration Using ECMWF and GFS Ensemble Reforecasts. Part I: Two-Meter

788        Temperatures. *Monthly Weather Review*, **136,** 2608-2619.

789    Hamill, T. M., R. Hagedorn, and J. S. Whitaker, 2008: Probabilistic Forecast

790        Calibration Using ECMWF and GFS Ensemble Reforecasts. Part II:

791        Precipitation. *Monthly Weather Review*, **136,** 2620-2632.

792    Hamill, T. M., and J. Juras, 2006: Measuring forecast skill: is it real skill or is it the

793        varying climatology? *Quarterly Journal of the Royal Meteorological Society*,

794        **132,** 2905-2923.

795    Hamill, T. M., and J. S. Whitaker, 2006: Probabilistic Quantitative Precipitation

796        Forecasts Based on Reforecast Analogs: Theory and Application. *Monthly*

797        *Weather Review*, **134,** 3209-3229.

798    ——, 2007: Ensemble Calibration of 500-hPa Geopotential Height and 850-hPa and

799        2-m Temperatures Using Reforecasts. *Monthly Weather Review*, **135,** 3273-

800        3280.

801    Hamill, T. M., J. S. Whitaker, M. Fiorino, and S. G. Benjamin, 2011: Global Ensemble

802        Predictions of 2009's Tropical Cyclones Initialized with an Ensemble Kalman

803        Filter. *Monthly Weather Review*, **139,** 668-688.

804    Hamill, T. M., J. S. Whitaker, and S. L. Mullen, 2006: Reforecasts: An Important

805        Dataset for Improving Weather Predictions. *Bulletin of the American*

806        *Meteorological Society*, **87,** 33-46.

807 Hamill, T. M., J. S. Whitaker, and X. Wei, 2004: Ensemble Reforecasting: Improving

808      Medium-Range Forecast Skill Using Retrospective Forecasts. *Monthly*

809      *Weather Review*, **132,** 1434-1447.

810 Higgins, R. W., J. E. Janowiak, and Y.-P. Yao, 1996: A Gridded Hourly Precipitation

811      Data Base for the United States (1963-1993). *NCEP/Climate Prediction Center*

812      *ATLAS No. 1, U. S. DEPARTMENT OF COMMERCE, National Oceanic and*

813      *Atmospheric Administration, National Weather Service.*

814 Hou, D., Z. Toth, Y. Zhu, and Y. Yang, 2008: Impact of a Stochastic Perturbation

815      Scheme on NCEP Global Ensemble Forecast System. *Proceedings, 19th AMS*

816      *conference on Probability and Statistics. New Orleans, LA, 20-24 Jan. 2008.*

817 Houtekamer, P. L., and H. L. Mitchell, 2005: Ensemble Kalman filtering. *Quarterly*

818      *Journal of the Royal Meteorological Society*, **131,** 3269-3289.

819 Houtekamer, P. L., H. L. Mitchell, and X. Deng, 2009: Model Error Representation in

820      an Operational Ensemble Kalman Filter. *Monthly Weather Review*, **137,** 2126-

821      2143.

822 Hunt, B., E. Kostelich, and I. Szunyogh, 2006: Efficient Data Assimilation for

823      Spatiotemporal Chaos: a Local Ensemble Transform Kalman Filter.

824 Isaksen, L., M. Bonavita, R. Buizza, M. Fisher, J. Haseler, M. Leutbecher, L. Raynaud,

825      2010: Ensemble of data assimilations at ECMWF, 48 pp.

826 Iversen, T., A. Deckmyn, C. Santos, K. A. I. Sattler, J. B. Bremnes, H. Feddersen, and I.-L.

827      Frogner, 2011: Evaluation of 'GLAMEPS'—a proposed multimodel EPS for

828      short range forecasting. *Tellus A*, **63,** 513-530.

829 Johnson, C., and R. Swinbank, 2009: Medium-range multimodel ensemble

830      combination and calibration. *Quarterly Journal of the Royal Meteorological*

831      *Society*, **135,** 777-794.

832 Kleist, D. T., D. F. Parrish, J. C. Derber, R. Treadon, W.-S. Wu, and S. Lord, 2009:

833      Introduction of the GSI into the NCEP Global Data Assimilation System.

834      *Weather and Forecasting*, **24,** 1691-1705.

835 Krishnamurti, T. N., and Coauthors, 2000: Multimodel Ensemble Forecasts for

836      Weather and Seasonal Climate. *Journal of Climate*, **13,** 4196-4216.

837 Leutbecher, M., 2005: On Ensemble Prediction Using Singular Vectors Started from

838      Forecasts. *Monthly Weather Review*, **133,** 3038-3046.

839 Lin, Y., and K. E. Mitchell, 2005: The NCEP Stage II/IV hourly precipitation analyses:

840      development and applications. *Preprints, 19th Conf. on Hydrology, American*

841      *Meteorological Society, San Diego, CA 9-13 January 2005, Paper 1.2.* .

842 Lorenz, J., H. Rauhut, F. Schweitzer, and D. Helbing, 2011: How social influence can

843      undermine the wisdom of crowd effect. *Proceedings of the National Academy*

844      *of Sciences*, **108,** 920-925.

845 Lowry, D. A., and H. R. Glahn, 1976: An Operational Model for Forecasting

846      Probability of Precipitation - PEATMOS PoP. *Monthly Weather Review*, **104,**

847      221-232.

848 Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF Ensemble

849      Prediction System: Methodology and validation. *Quarterly Journal of the*

850      *Royal Meteorological Society*, **122,** 73-119.

851    Mylne, K. R., R. E. Evans, and R. T. Clark, 2002: Multi-model multi-analysis ensembles

852        in quasi-operational medium-range forecasting. *Quarterly Journal of the*

853        *Royal Meteorological Society*, **128,** 361-384.

854    Palmer, T. N., and Coauthors, 2009: Stochastic parameterization and model

855        uncertainty. *ECMWF Tech Memo 589*.

856    Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski, 2005: Using Bayesian

857        Model Averaging to Calibrate Forecast Ensembles. *Monthly Weather Review*,

858        **133,** 1155-1174.

859    Rawlins, F., and Coauthors, 2007: The Met Office global four-dimensional variational

860        data assimilation scheme. *Quarterly Journal of the Royal Meteorological*

861        *Society*, **133,** 347-362.

862    Roulin, E., and S. Vannitsem, 2011: Post-processing of ensemble precipitation

863        predictions with extended logistic regression based on hindcasts. *Monthly*

864        *Weather Review*, **139,** Available from Emmanuel.Roulin@meteo.be.

865    Schmeits, M. J., and K. J. Kok, 2010: A Comparison between Raw Ensemble Output,

866        (Modified) Bayesian Model Averaging, and Extended Logistic Regression

867        Using ECMWF Ensemble Precipitation Reforecasts. *Monthly Weather Review*,

868        **138,** 4199-4211.

869    Shutts, G., 2005: A kinetic energy backscatter algorithm for use in ensemble

870        prediction systems. *Quarterly Journal of the Royal Meteorological Society*, **131,**

871        3079-3102.

872    Sloughter, J. M. L., A. E. Raftery, T. Gneiting, and C. Fraley, 2007: Probabilistic

873         Quantitative Precipitation Forecasting Using Bayesian Model Averaging.

874         *Monthly Weather Review*, **135,** 3209-3220.

875    Stensrud, D. J., and N. Yussouf, 2003: Short-Range Ensemble Predictions of 2-m

876         Temperature and Dewpoint Temperature over New England. *Monthly*

877         *Weather Review*, **131,** 2510-2524.

878    Tennant, W. J., G. J. Shutts, A. Arribas, and S. A. Thompson, 2011: Using a Stochastic

879         Kinetic Energy Backscatter Scheme to Improve MOGREPS Probabilistic

880         Forecast Skill. *Monthly Weather Review*, **139,** 1190-1206.

881    Vannitsem, S., and C. Nicolis, 2008: Dynamical Properties of Model Output Statistics

882         Forecasts. *Monthly Weather Review*, **136,** 405-419.

883    Vislocky, R. L., and J. M. Fritsch, 1995: Improved Model Output Statistics Forecasts

884         through Model Consensus. *Bulletin of the American Meteorological Society*, **76,**

885         1157-1164.

886    Wandishin, M. S., S. L. Mullen, D. J. Stensrud, and H. E. Brooks, 2001: Evaluation of a

887         Short-Range Multimodel Ensemble System. *Monthly Weather Review*, **129,**

888         729-747.

889    Wei, M., Z. Toth, R. Wobus, and Y. Zhu, 2008: Initial perturbations based on the

890         ensemble transform (ET) technique in the NCEP global operational forecast

891         system. *Tellus A*, **60,** 62-79.

892    Weigel, A. P., M. A. Liniger, and C. Appenzeller, 2008: Can multi-model combination

893         really enhance the prediction skill of probabilistic ensemble forecasts?

894         *Quarterly Journal of the Royal Meteorological Society*, **134,** 241-260.

895    Whitaker, J. S., X. Wei, and F. Vitart, 2006: Improving Week-2 Forecasts with

896          Multimodel Reforecast Ensembles. *Monthly Weather Review*, **134,** 2279-2284.

897    Wilks, D. S., 2006: *Statistical Methods in the Atmospheric Sciences (2nd Ed.).*

898          Academic Press, 627. pp.

899    ——, 2009: Extending logistic regression to provide full-probability-distribution

900          MOS forecasts. *Meteorological Applications*, **16,** 361-368.

901    Wilks, D. S., and T. M. Hamill, 2007: Comparison of Ensemble-MOS Methods Using

902          GFS Reforecasts. *Monthly Weather Review*, **135,** 2379-2390.

903    Wilson, L. J., S. Beauregard, A. E. Raftery, and R. Verret, 2007: Calibrated Surface

904          Temperature Forecasts from the Canadian Ensemble Prediction System

905          Using Bayesian Model Averaging. *Monthly Weather Review*, **135,** 1364-1385.

906    Yussouf, N., and D. J. Stensrud, 2007: Bias-Corrected Short-Range Ensemble

907          Forecasts of Near-Surface Variables during the 2005/06 Cool Season.

908          *Weather and Forecasting*, **22,** 1274-1286.

909

910
911

912 **Figure captions**
913

914
915 **Figure 1**: Illustration of the process for determining precipitation classes used in

916 the calculation of *BSS*. (a) Climatological probability of > 1-mm 24h$^{-1}$

917 precipitation as determined from Stage-IV data for September 2002-2009.

918 (b) Climatological class assigned to each grid point for September, 1-mm (24

919 h)$^{-1}$ event.

920 **Figure 2**: Brier skill scores of various forecasts for the (a) > 1-mm (24 h)$^{-1}$ event,

921 (b) > 10-mm (24 h)$^{-1}$ event, and (c) continuous ranked probability skill

922 scores, all as a function of forecast lead time. Error bars denote confidence

923 intervals, the 5$^{th}$ and 95$^{th}$ percentiles of a paired block bootstrap between

924 ECMWF and NCEP forecasts.

925 **Figure 3**: Maps of average *CRPSS* for day +3 forecasts for (a) ECMWF, (b) NCEP, (c)

926 UKMO, and (d) CMC.

927 **Figure 4**: (a) RMS errors, and (b) bias for day +3 forecasts, each as a function of the

928 climatological probability of greater than 1-mm (24 h)$^{-1}$. Light grey bars in

929 panel (a) denote the relative frequency of each climatological probability.

930 **Figure 5**: Reliability diagrams for day +3 forecasts for the > 10-mm (24 h)$^{-1}$ event.

931 (a) ECMWF, (b) NCEP, (c) CMC, and (d) UKMO. The dark line on each is the

932 20-member reliability curve. The lighter grey line on panel (a) is the

933 reliability for the full 50-member ensemble. The inset histogram bars show

934 the relative frequency of usage for each probability bin. The black lines on

935 the inset are the relative frequency of usage for the climatological

936        distribution across all the sample points.  The grey dots on the inset

937        histogram of panel (a) are the relative frequency of usage for the ECMWF full

938        50-member ensemble.

939    **Figure 6**: Analyzed  > 10-mm $(24\,h)^{-1}$ precipitation boundary (black line) and area

940        exceeding 10 mm (grey shading) for 25 cases with the largest areal coverage

941        of greater than 10 mm in the upper Midwest US.  Red lines indicate the 0.5

942        probability contour from the ECMWF ensemble for the day +3 forecasts of >

943        10 mm $(24\,h)^{-1}$.

944    **Figure 7**: (a) Analyzed precipitation for the 24-h period ending 00 UTC 21 July 2010.

945        10-mm $(24\,h)^{-1}$ contour is denoted by the thick black line.  (b) Probability of

946        greater than 10 mm $(24\,h)^{-1}$ for day +3 forecast from the ECMWF ensemble

947        for the same period.  The analyzed 10-mm contour from panel (a) is repeated.

948        (c) as in (b), but for NCEP.  (d) CMC, (e) UK Met Office, and (f) multi-model

949        combination.

950    **Figure 8**: As in Fig. 7, but for 24-h period ending 00 UTC 8 August 2010.

951    **Figure 9**:  Brier skill scores of various forecasts for (a) > 1-mm $(24\,h)^{-1}$ event, and

952        (b) > 10-mm $(24\,h)^{-1}$ event, and (c) continuous ranked probability skill

953        scores, all as a function of forecast lead time.  "Multi-model/cal" refers to

954        forecasts from the multi-model, calibrated using ELR.  "ECMWF/reforecast"

955        refers to ECMWF forecasts calibrated using ELR and the reforecast data set.

956        Error bars denote confidence intervals, the 5th and 95th percentiles of a

957        paired block bootstrap between ECMWF and NCEP forecasts.

958    **Figure 10**:  Maps of *CRPSS* for day +3 forecasts for (a) multi-model, (b) multi-model

959          with ELR calibration, and (c) ECMWF with ELR calibration using reforecasts.

960    **Figure 11**: As in Fig. 3, reliability diagrams for day +3 forecasts at the > 10-mm

961          $(24 \text{ h})^{-1}$ event. (a) Multi-model forecasts, (b) multi-model with ELR

962          calibration, and (c) ECMWF with ELR calibration using reforecasts.

963    **Figure 12**:  A histogram of the absolute errors of day +3 ensemble-mean

964          precipitation forecasts for the 2002 and 2006 reforecasts and for the 2010,

965          20-member real-time ensemble.

966    **Figure 13**:  As in Fig. 6, but for multi-model forecasts.
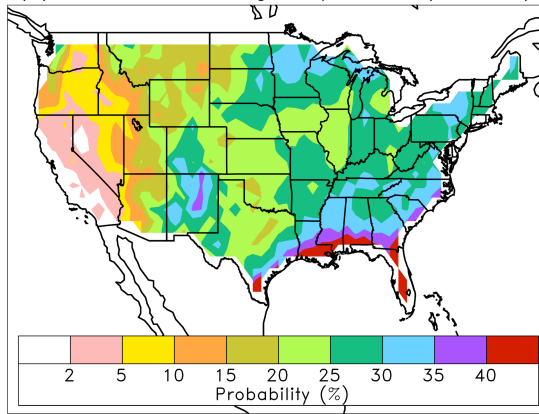
967    **Figure 14**:  As in Fig. 6, but for reforecast-calibrated ECMWF forecasts.

968    **Figure 15**:  (a) Analyzed precipitation for the 24-h period ending 00 UTC 21 July

969          2010. 10-mm contour is denoted by the thick black line.  (b) Probability of

970          greater than 10 mm $(24 \text{ h})^{-1}$ for day +3 forecast from the ECMWF ensemble

971          for the same period.  (c) as in (b), but for multi-model ensemble, and (d) as in

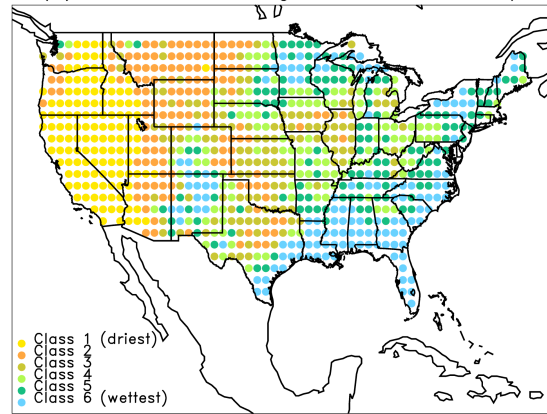972          (b), but for reforecast-calibrated ECMWF ensemble.

973    **Figure 16**:  As in Fig. 15, but for the 24-h period ending 00 UTC 8 August 2010.

974
975
976
977
978
979
980

981
982
983



(a) 1−mm climatological probability for Sep
(b) 1−mm climatological classes for Sep

Probability (%)
2  5  10  15  20  25  30  35  40

Class 1 (driest)
Class 2
Class 3
Class 4
Class 5
Class 6 (wettest)

984
985
986  **Figure 1**:  Illustration of the process for determining precipitation classes used in
987  the calculation of *BSS*.  (a) Climatological probability of > 1-mm 24h$^{-1}$ precipitation
988  as determined from Stage-IV data for September 2002-2009.  (b) Climatological
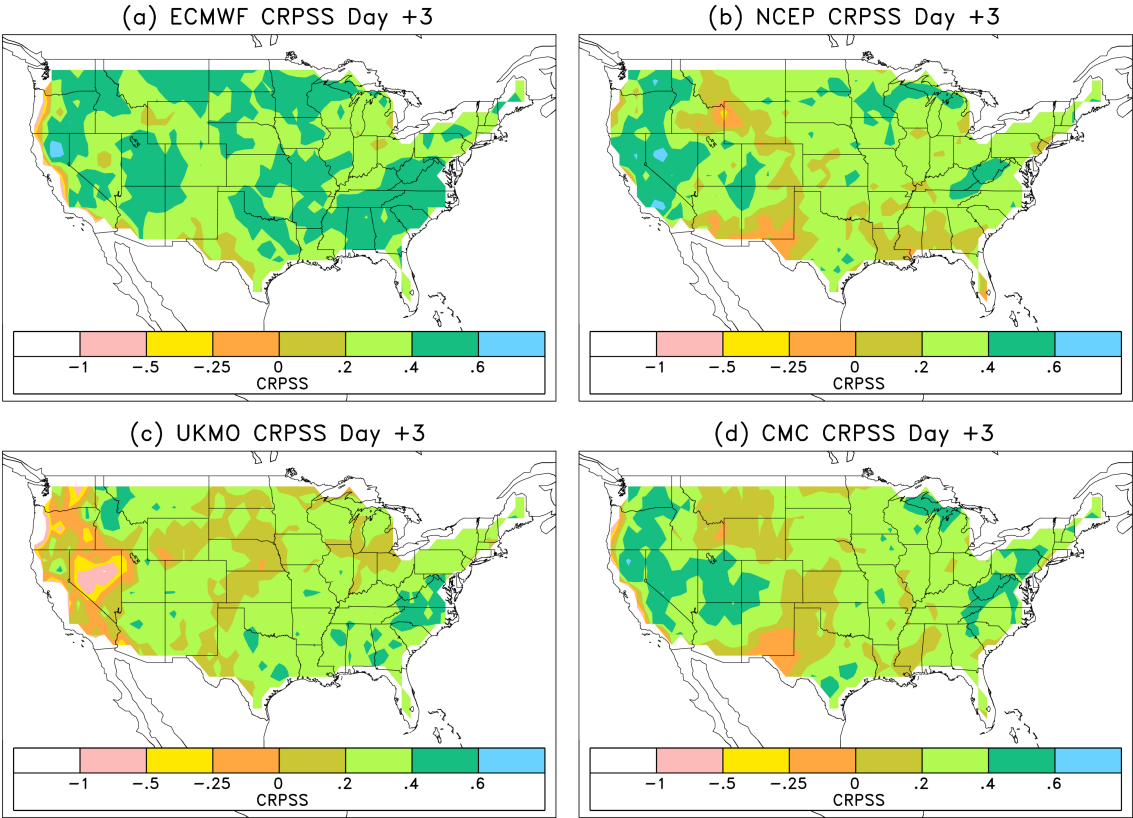989  class assigned to each grid point for September, 1-mm (24 h)$^{-1}$ event.
990
991

992
993



(a) Brier Skill Scores, 1 mm

(b) Brier Skill Scores, 10 mm

ECMWF
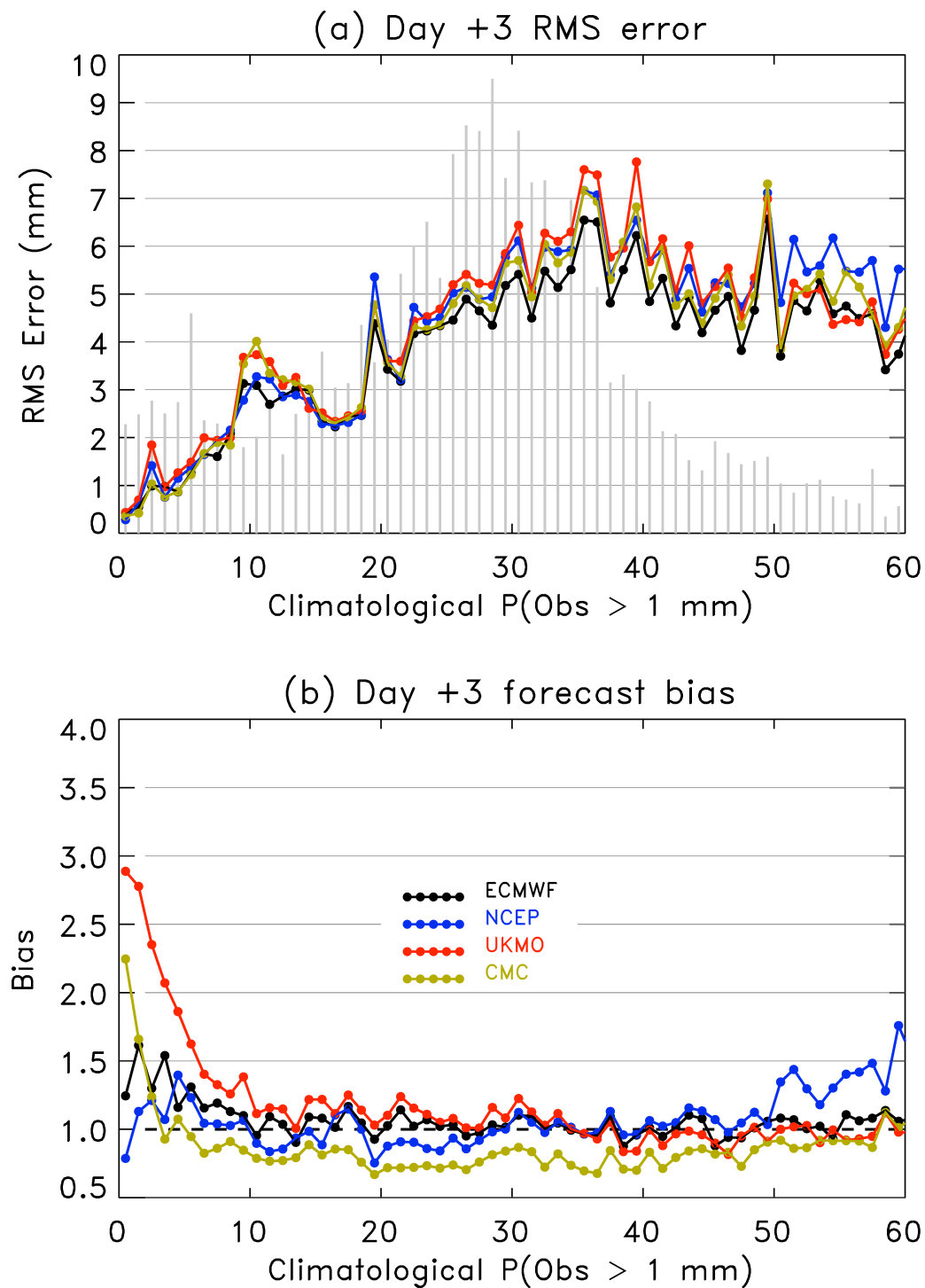NCEP
UKMO
CMC

(c) Continuous Ranked Probability Skill Scores

994
995 **Figure 2**: Brier skill scores of various forecasts for the (a) > 1-mm (24 h)$^{-1}$ event,
996 (b) > 10-mm (24 h)$^{-1}$ event, and (c) continuous ranked probability skill scores, all as
997 a function of forecast lead time. Error bars denote confidence intervals, the 5[th] and
998 95[th] percentiles of a paired block bootstrap between ECMWF and NCEP forecasts.
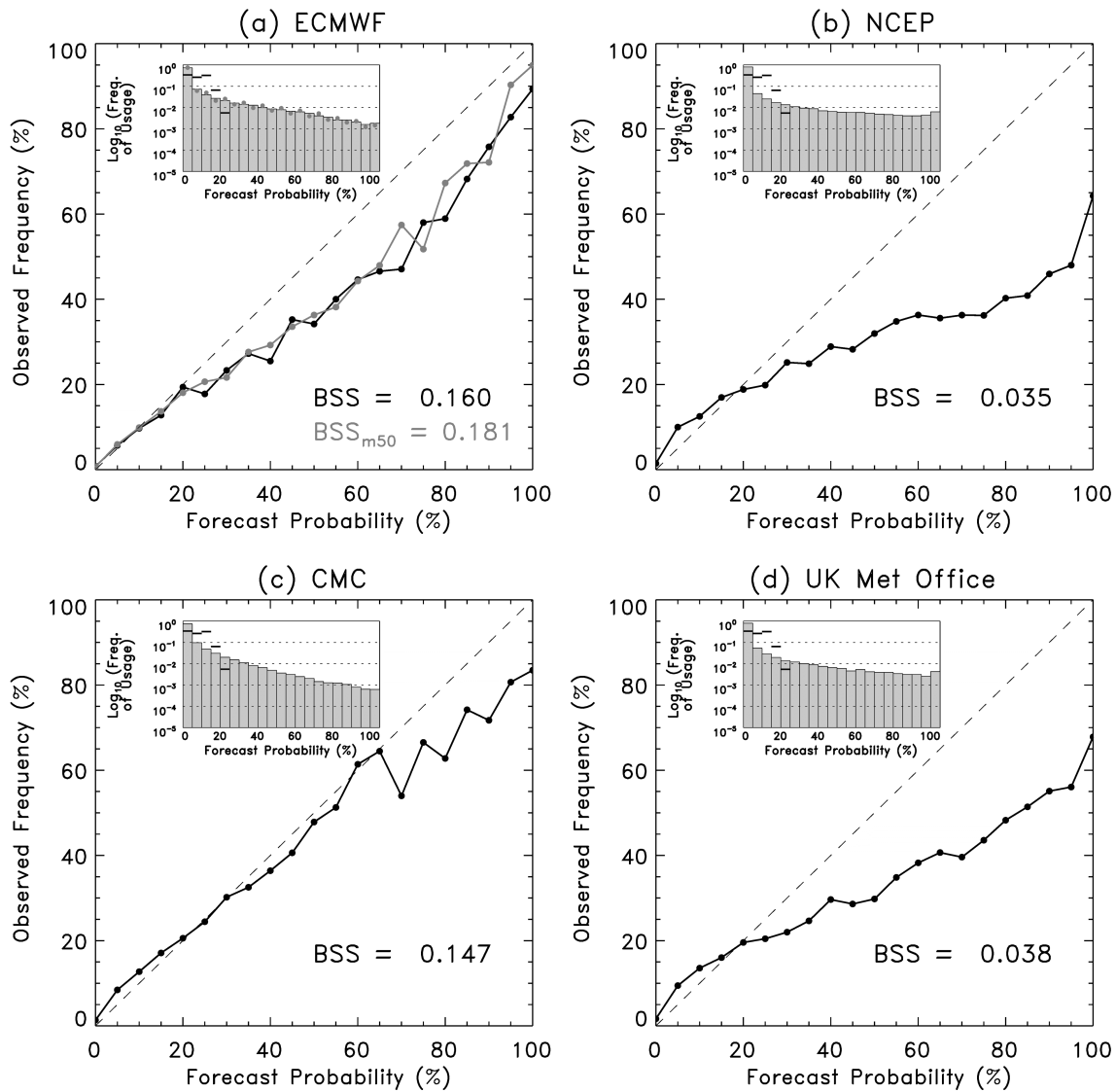999

1000

(a) ECMWF CRPSS Day +3    (b) NCEP CRPSS Day +3

(c) UKMO CRPSS Day +3    (d) CMC CRPSS Day +3

1001
1002
1003    **Figure 3**: Maps of average *CRPSS* for day +3 forecasts for (a) ECMWF, (b) NCEP, (c)
1004    UKMO, and (d) CMC.
1005

## (a) Day +3 RMS error
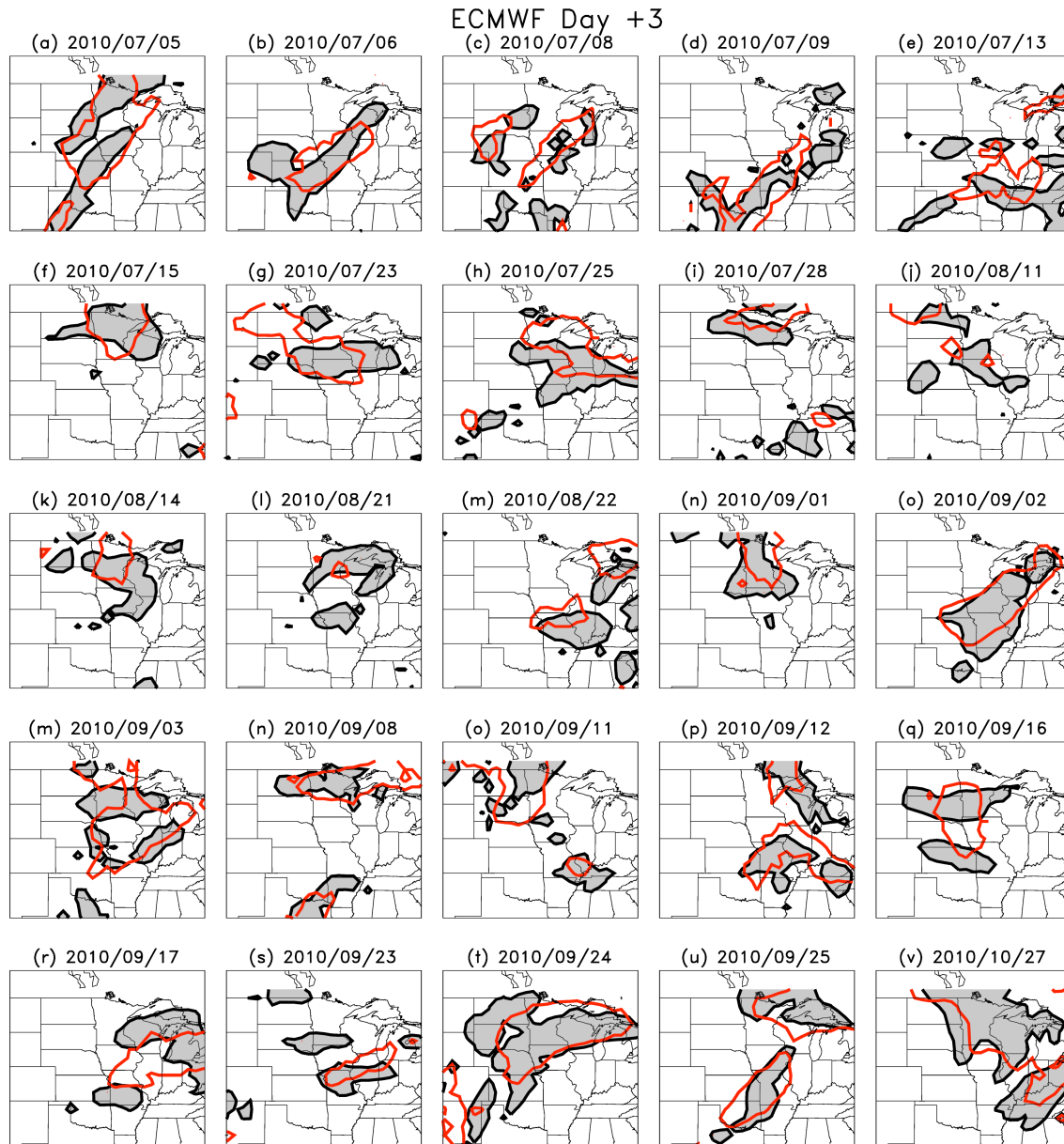
## (b) Day +3 forecast bias

1010  **Figure 4**:  (a) RMS errors, and (b) bias for day +3 forecasts, each as a function of the
1011  climatological probability of greater than 1-mm (24 h)$^{-1}$.  Light grey bars in panel (a)
1012  denote the relative frequency of each climatological probability.

Figure 5: Reliability diagrams for day +3 forecasts for the > 10-mm (24 h)$^{-1}$ event.
(a) ECMWF, (b) NCEP, (c) CMC, and (d) UKMO.  The dark line on each is the 20-
member reliability curve.  The lighter grey line on panel (a) is the reliability for the
full 50-member ensemble.  The inset histogram bars show the relative frequency of
usage for each probability bin.  The black lines on the inset are the relative
frequency of usage for the climatological distribution across all the sample
points.  The grey dots on the inset histogram of panel (a) are the relative frequency
of usage for the ECMWF full 50-member ensemble.

1025
1026

ECMWF Day +3

(a) 2010/07/05  (b) 2010/07/06  (c) 2010/07/08  (d) 2010/07/09  (e) 2010/07/13

(f) 2010/07/15  (g) 2010/07/23  (h) 2010/07/25  (i) 2010/07/28  (j) 2010/08/11

(k) 2010/08/14  (l) 2010/08/21  (m) 2010/08/22  (n) 2010/09/01  (o) 2010/09/02

(m) 2010/09/03  (n) 2010/09/08  (o) 2010/09/11  (p) 2010/09/12  (q) 2010/09/16

(r) 2010/09/17  (s) 2010/09/23  (t) 2010/09/24  (u) 2010/09/25  (v) 2010/10/27
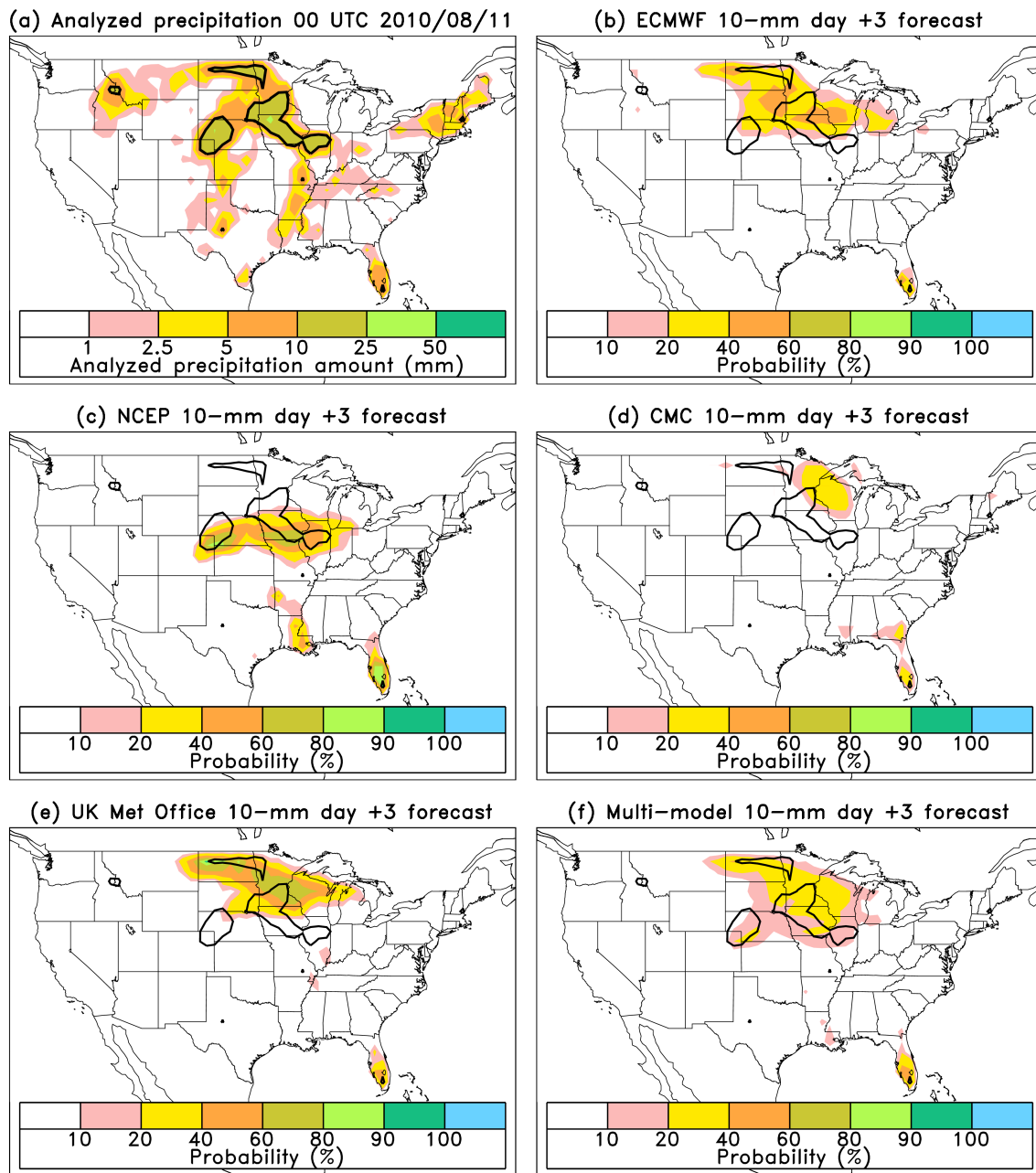
1027
1028

1029  **Figure 6**: Analyzed > 10-mm (24 h)$^{-1}$ precipitation boundary (black line) and area
1030  exceeding 10 mm (grey shading) for 25 cases with the largest areal coverage of
1031  greater than 10 mm in the upper Midwest US.  Red lines indicate the 0.5 probability
1032  contour from the ECMWF ensemble for the day +3 forecasts of > 10 mm (24 h)$^{-1}$.
1033

1034



(a) Analyzed precipitation 00 UTC 2010/07/21
(b) ECMWF 10−mm day +3 forecast
(c) NCEP 10−mm day +3 forecast
(d) CMC 10−mm day +3 forecast
(e) UK Met Office 10−mm day +3 forecast
(f) Multi−model 10−mm day +3 forecast

1035
1036 **Figure 7**: (a) Analyzed precipitation for the 24-h period ending 00 UTC 21 July 2010.
1037 10-mm $(24 \text{ h})^{-1}$ contour is denoted by the thick black line. (b) Probability of greater
1038 than 10 mm $(24 \text{ h})^{-1}$ for day +3 forecast from the ECMWF ensemble for the same
1039 period. The analyzed 10-mm contour from panel (a) is repeated. (c) as in (b), but
1040 for NCEP. (d) CMC, (e) UK Met Office, and (f) multi-model combination.
1041

51

1042



(a) Analyzed precipitation 00 UTC 2010/08/11

(b) ECMWF 10−mm day +3 forecast

(c) NCEP 10−mm day +3 forecast

(d) CMC 10−mm day +3 forecast

(e) UK Met Office 10−mm day +3 forecast

(f) Multi−model 10−mm day +3 forecast

1043
1044
1045 **Figure 8**: As in Fig. 7, but for 24-h period ending 00 UTC 8 August 2010.
1046

1047
1048



(a) Brier Skill Scores, 1 mm

(b) Brier Skill Scores, 10 mm

Multi-model / Cal
Multi-model
ECMWF
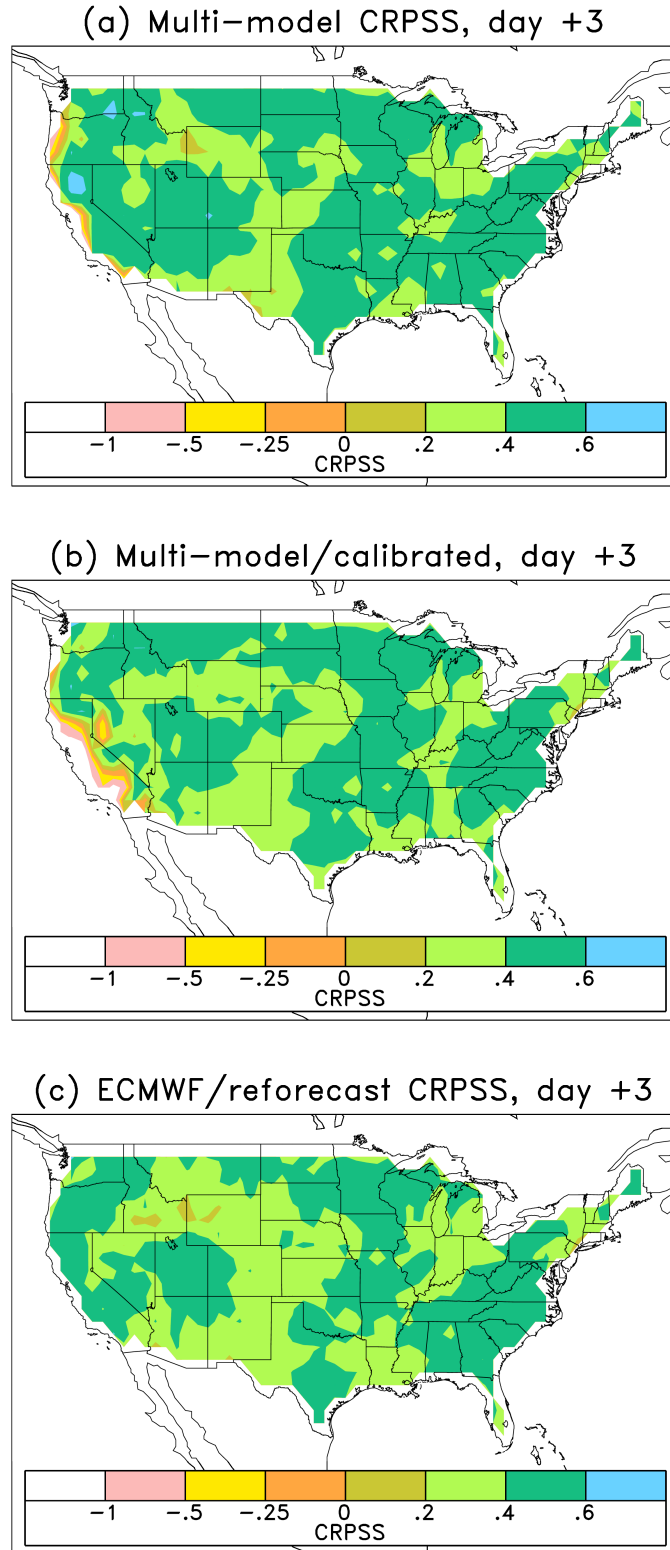ECMWF / reforecast
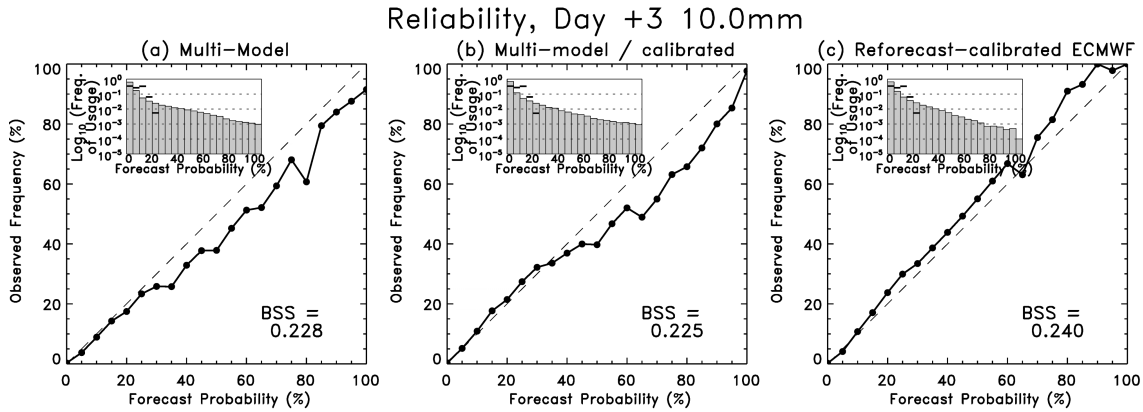
(c) Continuous Ranked Probability Skill Scores

1049
1050   **Figure 9**:  Brier skill scores of various forecasts for (a) > 1-mm (24 h)$^{-1}$ event, and
1051   (b) > 10-mm (24 h)$^{-1}$ event, and (c) continuous ranked probability skill scores, all as
1052   a function of forecast lead time.  "Multi-model/cal" refers to forecasts from the
1053   multi-model, calibrated using ELR.  "ECMWF/reforecast" refers to ECMWF forecasts
1054   calibrated using ELR and the reforecast data set. Error bars denote confidence
1055   intervals, the 5$^{th}$ and 95$^{th}$ percentiles of a paired block bootstrap between ECMWF
1056   and NCEP forecasts.
1057

## (a) Multi-model CRPSS, day +3

## (b) Multi-model/calibrated, day +3

## (c) ECMWF/reforecast CRPSS, day +3
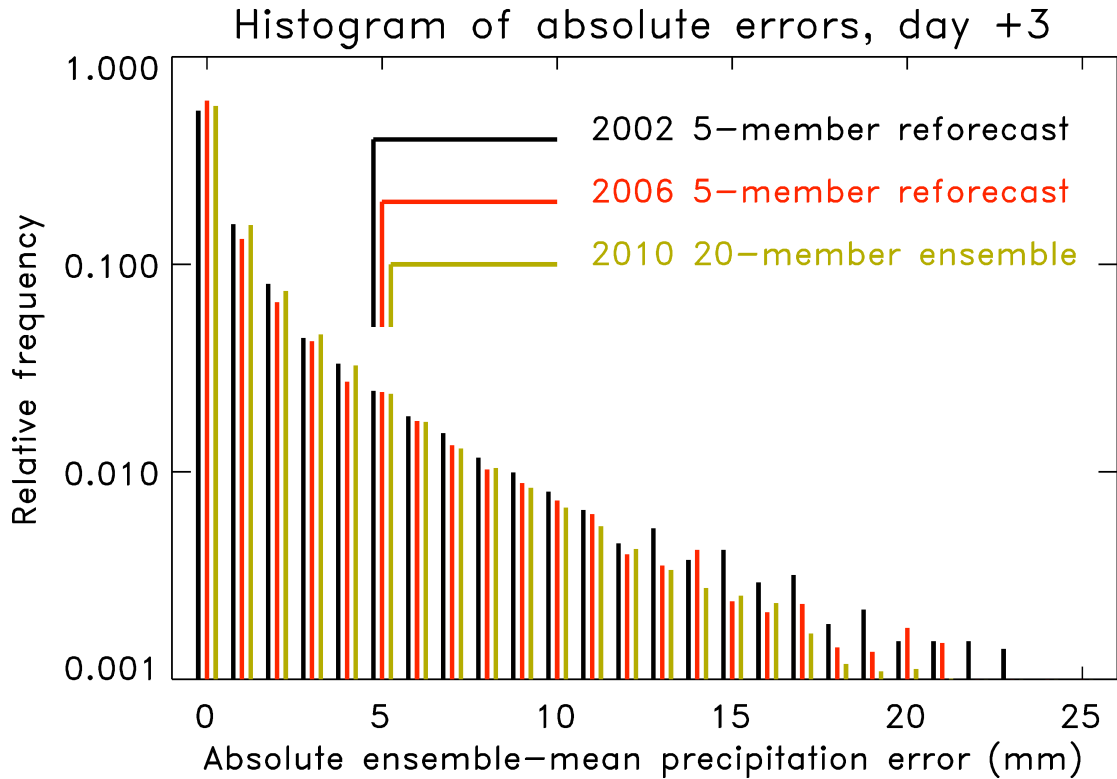
1058
1059
1060   **Figure 10**: Maps of *CRPSS* for day +3 forecasts for (a) multi-model, (b) multi-model
1061   with ELR calibration, and (c) ECMWF with ELR calibration using reforecasts.

Reliability, Day +3 10.0mm

1062
1063
1064 **Figure 11**: As in Fig. 3, reliability diagrams for day +3 forecasts at the > 10-mm
1065 (24 h)$^{-1}$ event. (a) Multi-model forecasts, (b) multi-model with ELR calibration, and
1066 (c) ECMWF with ELR calibration using reforecasts.
1067
1068



Histogram of absolute errors, day +3

1069
1070
1071 **Figure 12**:  A histogram of the absolute errors of day +3 ensemble-mean
1072 precipitation forecasts for the 2002 and 2006 reforecasts and for the 2010, 20-
1073 member real-time ensemble.
1074
1075

1076

Multi−Model Day +3

(a) 2010/07/05   (b) 2010/07/06   (c) 2010/07/08   (d) 2010/07/09   (e) 2010/07/13

(f) 2010/07/15   (g) 2010/07/23   (h) 2010/07/25   (i) 2010/07/28   (j) 2010/08/11

(k) 2010/08/14   (l) 2010/08/21   (m) 2010/08/22   (n) 2010/09/01   (o) 2010/09/02

(m) 2010/09/03   (n) 2010/09/08   (o) 2010/09/11   (p) 2010/09/12   (q) 2010/09/16

(r) 2010/09/17   (s) 2010/09/23   (t) 2010/09/24   (u) 2010/09/25   (v) 2010/10/27

1077
1078
1079   **Figure 13**:  As in Fig. 6, but for multi-model forecasts.
1080

56

1081

ECMWF/reforecast Day +3



**Figure 14**:  As in Fig. 6, but for reforecast-calibrated ECMWF forecasts.

1082
1083
1084
1085

57

1086



(a) Analyzed precipitation, 00 UTC 2010/07/21

(b) ECMWF 10-mm day +3 forecast

Analyzed precipitation amount (mm)
1    2.5    5    10    25    50

Probability (%)
10    20    40    60    80    90    100

(c) Multi-model 10-mm day +3 forecast

(d) ECMWF/reforecast 10-mm day +3 forecast

Probability (%)
10    20    40    60    80    90    100

Probability (%)
10    20    40    60    80    90    100

1087
1088
1089  **Figure 15**: (a) Analyzed precipitation for the 24-h period ending 00 UTC 21 July
1090  2010. 10-mm contour is denoted by the thick black line. (b) Probability of greater
1091  than 10 mm $(24 \text{ h})^{-1}$ for day +3 forecast from the ECMWF ensemble for the same
1092  period. (c) as in (b), but for multi-model ensemble, and (d) as in (b), but for
1093  reforecast-calibrated ECMWF ensemble.
1094
1095

58

1096



(a) Analyzed precipitation, 00 UTC 2010/08/11

(b) ECMWF 10−mm day +3 forecast

Analyzed precipitation amount (mm)

1    2.5    5    10    25    50

Probability (%)

10    20    40    60    80    90    100

(c) Multi−model 10−mm day +3 forecast

(d) ECMWF/reforecast 10−mm day +3 forecast

Probability (%)

10    20    40    60    80    90    100

Probability (%)

10    20    40    60    80    90    100

1097
1098
1099    **Figure 16**:  As in Fig. 15, but for the 24-h period ending 00 UTC 8 August 2010.
1100